

Thèse de Doctorat

Moustafa IBRAHIM

*Mémoire présenté en vue de l'obtention du
grade de Docteur de l'École Centrale de Nantes
sous le label de l'Université de Nantes Angers Le Mans*

École doctorale : Sciences et Technologies de l'Information et Mathématiques

Discipline : Mathématiques et leurs interactions

Unité de recherche : Laboratoire de Mathématiques Jean Leray (LMJL)

Soutenue le 25 septembre 2014

Systèmes paraboliques dégénérés intervenant en mécanique des fluides et en médecine: analyse mathématique et numérique

JURY

Président :	M. Iraj MORTAZAVI , Professeur, Conservatoire National des Arts et Métiers
Rapporteurs :	M. Laurent CHUPIN , Professeur, Université Blaise Pascal M. Robert EYMARD , Professeur, Université Paris-Est Marne-la-Vallée
Examineurs :	M. Christophe BERTHON , Professeur, Université de Nantes M. Laurent DI MENZA , Professeur, Université de Reims
Invité :	M. Clément CANCÈS , Maître de conférences, Université Paris VI
Directeur de thèse :	M. Mazen SAAD , Professeur, École Centrale de Nantes

Remerciements

Il m'a été donné de rencontrer une personne qui allait être pour moi une présence remarquable par sa richesse d'humanité. C'est avec une telle qualité de présence, toujours bienveillante et attentive, que Mazen Saad a donné réalité à mon désir d'explorer un passionnant domaine entre les mathématiques et la médecine. Il a toujours fait preuve d'une totale confiance, d'un soutien infiniment précieux et d'une gentillesse permanente à mon égard, qui m'ont permis de mener et d'achever mon travail de recherche pendant ces trois années. Je suis très heureux d'avoir eu la chance de bénéficier au niveau scientifique et humain de sa grande expérience, de sa vision mathématique, ainsi que de son immense culture. Trouver des mots assez forts pour dire "Merci" est une chose difficile, d'autant plus dans une langue qui n'est pas ma langue maternelle. Je souhaite ici remercier Mazen Saad pour son amitié et lui exprimer ma plus profonde gratitude.

Je souhaite remercier très vivement les professeurs qui se sont intéressés de près à mon travail : Laurent Chupin et Robert Eymard, qui ont accepté d'en être les rapporteurs et de donner de leur temps pour examiner cette thèse. Je remercie également Christophe Berthon, Laurent Di Menza et Iraj Mortazavi qui m'ont tous fait l'honneur et le plaisir de constituer le jury de ma thèse.

J'adresse également de vifs remerciements à Clément Cancès pour avoir accepté l'invitation et participer au jury de ma thèse, ainsi que pour sa contribution dans le quatrième chapitre de la thèse. Je le remercie pour son amitié et pour les nombreuses discussions que nous avons eues lors de sa visite à l'École Centrale de Nantes, au cours de laquelle il n'a pas hésité à partager avec moi son savoir et son expérience dans les domaines numérique et théorique. Je lui en suis très reconnaissant.

Je voudrais également exprimer mes plus chaleureux remerciements à Françoise Foucher, pour son aide et ses conseils pendant les trois années de la thèse. Je la remercie pour m'avoir donné l'occasion d'enseigner à l'École Centrale de Nantes, ainsi je n'oublie pas ses remarques et sa contribution à corriger les fautes d'orthographe dans l'introduction de la thèse.

Je remercie de façon générale les membres du laboratoire de mathématiques Jean Leray ainsi que les membres du département informatique et mathématiques de l'École Centrale de Nantes, dans lequel j'ai été accueilli avec bonne humeur et convivialité.

Au delà des horizons nantais, je tiens également à saluer et remercier toutes les personnes rencontrées dans les divers congrès ou séjours de recherche avec qui il a été agréable d'échanger ou de travailler. Je souhaite ici remercier Georges Chamoun pour son amitié et pour les beaux moments qu'on a passé ensemble lors de notre participation aux congrès à Sofia, Marseille et Taormina.

M'éloignant de la sphère des mathématiques, j'adresse ma reconnaissance et bien plus... à tous mes amis (en France et au Liban) pour leur soutien. Sans eux je n'aurais pas eu la joie de conti-

nuer sur ce long chemin studieux. Ce “succès” est également un peu le leur. Un grand Merci à Rita-Maria, Hanna, Rabih, Mohamad, Chady, Georges Salameh, Georges Massaad, Ali, Ibrahim, Michaël, Olivier... pour leur présence à ma soutenance et pour leur immense aide dans l’organisation du pot de la soutenance. Et je n’oublie pas à saluer Amine, Akram, Salim, Tala, Youssef, Hala,... pour leurs encouragements. À eux tous, je souhaite beaucoup de bien.

Je clos enfin ces remerciements en dédiant cette thèse de doctorat à ma famille : mon père, ma mère, mes frères (Ibrahim, Wissam, Abed, Majdi et Maher). Grâce à eux mon séjour en France et la poursuite de mes études jusqu’à ce stade ont été possibles. Je les remercie de tout mon cœur pour leur amour et leur soutien discret et essentiel. C’est avec leurs encouragements que j’ai pu accomplir ce travail de thèse. Il m’est impossible de trouver des mots pour dire à quel point je suis fier d’eux, et à quel point je les aime. J’en profite pour glisser un clin d’œil aussi au reste de la famille (au Liban, Syrie et Australie), sans qui je n’aurais pas eu le goût de me lancer dans cette aventure.

“At each given moment, there is only a fine layer between the “trivial” and the impossible. Mathematical discoveries are made in this layer.”

A.N. Kolmogorov

Résumé

Dans cette thèse, nous nous intéressons à l'analyse mathématique et numérique des systèmes paraboliques non linéaires dégénérés découlant, soit de la modélisation de la chimiotaxie, soit de la modélisation des fluides compressibles. Le modèle de chimiotaxie (Keller-Segel) proposé est un modèle de dynamique des populations décrivant l'évolution spatio-temporelle de la densité cellulaire et de la concentration chimiotactique. Pour ce modèle, nous étudions la formation de patterns en utilisant l'analyse de stabilité linéaire et le principe de Turing. Nous proposons ensuite un schéma numérique CVFE pour un modèle anisotrope de Keller-Segel. La construction de ce schéma est basée sur la méthode des éléments finis pour le terme de diffusion et sur la méthode des volumes finis classique pour le terme de convection. Nous montrons que ce schéma assure le principe de maximum discret et qu'il est consistant dans le cas où tous les coefficients de transmissibilité sont positifs. Par la suite, sur des maillages triangulaires généraux, nous proposons et analysons un schéma numérique CVFE non linéaire. Ce schéma est basé sur l'utilisation d'un flux numérique de Godunov pour le terme de diffusion, tandis que le terme de convection est approché au moyen d'un décentrage amont et d'un flux de Godunov. D'une part, le décentrage amont permet d'avoir le principe de maximum. D'autre part, le flux de Godunov assure que les solutions discrètes soient bornées sans restriction sur le maillage du domaine spatial ni sur les coefficients de transmissibilité. Nous réalisons différentes simulations numériques bi-dimensionnelles pour illustrer l'efficacité du schéma à tenir compte des hétérogénéités. Enfin, nous nous intéressons à une équation parabolique dégénérée contenant des termes dégénérés d'ordre 0 et 1 et décrivant un modèle de chimiotaxie-fluide ou l'écoulement d'un fluide compressible. Une formulation faible classique est souvent possible en absence des termes dégénérés d'ordre 0 et 1 ; tandis que dans le cas général, nous obtenons des solutions dans un sens affaibli vérifiant une formulation de type inégalité variationnelle. La définition des solutions faibles est adaptée à la nature de la dégénérescence des termes de dissipation.

Mots clés

Systèmes de réaction-diffusion, Formation de patterns, Chimiotaxie, Stabilité linéaire, Méthode de volumes finis, Méthode des éléments finis, Systèmes paraboliques dégénérés, Tenseurs anisotropes hétérogènes, Fluide compressible.

Abstract

In this thesis, we are interested in the mathematical and numerical analysis of nonlinear degenerate parabolic systems arising either from modeling the chemotaxis process, or from modeling compressible flows in porous media. The proposed chemotaxis model (Keller-Segel model) is a model of population dynamics describing the spatio-temporal evolution of the cell density and the chemical concentration. For this model, we study the pattern formation using the linear stability analysis as well as the principle of Turing. Then, we propose a numerical scheme (CVFE scheme) for an anisotropic Keller-Segel model. The construction of the scheme is based on the use of each of the finite element scheme for the diffusion term and the upwind finite volume scheme for the convective term. We show that the scheme is consistent and ensures the discrete maximum principle in the case where all the transmissibility coefficients are nonnegative. Thereafter, over general triangular meshes, we propose and analyze a nonlinear CVFE scheme. This scheme is based on the use of the Godunov flux function for the diffusion term, while the convective term is approximated by parts using an upwind finite volume scheme and a Godunov flux function. First, the upwind finite volume scheme allows of having the discrete maximum principle. On the other hand, the Godunov scheme ensures the boundedness of the discrete solutions without restrictions on the mesh nor on the transmissibility coefficients. Using this scheme, we realize some numerical simulations to illustrate the effectiveness of the scheme. Finally, we are interested in a degenerate parabolic equation containing degenerate terms of order 0 and 1 and describing a chemotaxis-fluid model or a displacement of compressible flows. Classical weak formulation is often possible in the absence of degenerate terms of order 0 and 1 ; while in the general case, we obtain weak solutions in the sense of verifying a weighted formulation. The definition of weak solutions is adapted to the nature of the degeneracy of the dissipative terms.

Keywords

Reaction-diffusion systems, Pattern formation, Chemotaxis, Linear stability, Finite volume method, Finite element method, Degenerate parabolic systems, Heterogeneous anisotropic tensors, Compressible fluid.

Table des matières

1	Introduction	1
1.1	Contexte général et scientifique	1
1.1.1	Motivation biologique	2
1.1.2	Motivation mathématique (Principe de Turing)	2
1.1.3	Validation de l'idée de Turing	3
1.1.4	La chimiotaxie	3
1.2	Plan de la thèse	4
1.2.1	Chapitre 2 : Formation de patterns pour un modèle de chimiotaxie	4
1.2.2	Chapitre 3 : Un schéma volumes finis éléments finis pour capturer les patterns pour un modèle de chimiotaxie	10
1.2.3	Chapitre 4 : Un schéma CVFE non linéaire pour le modèle de Keller–Segel modifié	17
1.2.4	Chapitre 5 : Analyse d'une équation parabolique dégénérée modélisant la chimiotaxie ou les fluides compressibles en milieu poreux	25
2	Pattern formation for a volume-filling chemotaxis model	29
2.1	Introduction	29
2.2	Volume-filling chemotaxis model	31
2.3	Pattern formation	33
2.3.1	Linear stability	33
2.3.2	Formal asymptotic expansion	34
2.3.3	Turing conditions	34
2.3.4	Bifurcation	36
2.4	Analysis for a nonlinear density	37
2.4.1	Bifurcation with chemotactic sensitivity χ	38
2.4.2	Bifurcation with growth rate α	40
2.4.3	Bifurcation with death rate β	41
2.5	Finite volume approximation	41
2.5.1	Space-time discretization and discrete functions	42
2.5.2	Finite volume scheme for system (2.29)	45
2.5.3	Numerical results	47
3	Capture of patterns for an anisotropic chemotaxis model	53
3.1	Introduction	53
3.2	Volume-filling chemotaxis model	55
3.3	CVFE discretization of the continuous problem	58
3.3.1	CVFE scheme for the modified Keller–Segel model	60
3.3.2	Main result	63
3.4	<i>A priori</i> analysis of discrete solutions	64

3.4.1	Nonnegativity of $v_{\mathcal{M},\Delta t}$, confinement of $u_{\mathcal{M},\Delta t}$	64
3.4.2	Discrete <i>a priori</i> estimates	66
3.4.3	Existence of a discrete solution	69
3.5	Compactness estimates on discrete solutions	71
3.5.1	Time translate estimate	71
3.5.2	Space translate estimate	73
3.6	Convergence of the CVFE scheme	74
3.7	Numerical simulation in two-dimensional space	79
4	A nonlinear CVFE scheme for the modified Keller-Segel model	85
4.1	Introduction	86
4.2	The modified chemotaxis model	86
4.3	Space-time discretization and notations	88
4.3.1	Space discretizations of Ω .	88
4.3.2	Discrete finite elements space $\mathcal{H}_{\mathcal{T}}$, control volumes space $\mathcal{X}_{\mathcal{M}}$.	88
4.3.3	Time discretization of $(0, t_f)$.	89
4.3.4	Space-time discretization of Q_{t_f} .	89
4.3.5	Main property	90
4.4	The nonlinear CVFE scheme	90
4.4.1	Discretization of the first equation of system (4.1)	91
4.4.2	Discretization of the second equation of system (4.1)	93
4.4.3	Main result	94
4.5	Discrete properties, a priori estimates and existence	94
4.5.1	Discrete maximum principle	97
4.5.2	Entropy estimates on $v_{\mathcal{M},\Delta t}$	99
4.5.3	Energy estimates on $u_{\mathcal{M},\Delta t}$	101
4.5.4	Enhanced estimate on $v_{\mathcal{M},\Delta t}$	103
4.5.5	Existence of a discrete solution	104
4.6	Compactness estimates on the family of discrete solutions.	106
4.6.1	Time translate estimate.	106
4.6.2	Space translate estimate.	108
4.7	Convergence	109
4.7.1	Identification as a weak solution	110
4.8	Numerical results	114
5	Study of a degenerate parabolic nonlinear equation	119
5.1	Introduction	119
5.2	The nonlinear degenerate model	120
5.2.1	Classical weak solutions	121
5.2.2	Weak degenerate solutions	122
5.3	Existence for the nondegenerate case	123
5.3.1	Weak nondegenerate solutions	124
5.3.2	Maximum principle on the saturation	128
5.4	Proof of theorem 5.2.	130
5.5	Proof of theorem 5.4	134
5.6	Proof of theorem 5.6	145
A	Technical Lemmas	151
A.1	The reference element	154

Bibliography**156**

Introduction

Sommaire

1.1	Contexte général et scientifique	1
1.1.1	Motivation biologique	2
1.1.2	Motivation mathématique (Principe de Turing)	2
1.1.3	Validation de l'idée de Turing	3
1.1.4	La chimiotaxie	3
1.2	Plan de la thèse	4
1.2.1	Chapitre 2 : Formation de patterns pour un modèle de chimiotaxie	4
1.2.2	Chapitre 3 : Un schéma volumes finis éléments finis pour capturer les patterns pour un modèle de chimiotaxie	10
1.2.3	Chapitre 4 : Un schéma CVFE non linéaire pour le modèle de Keller–Segel modifié	17
1.2.4	Chapitre 5 : Analyse d'une équation parabolique dégénérée modélisant la chimiotaxie ou les fluides compressibles en milieu poreux	25

1.1 Contexte général et scientifique

Les êtres humains ont longtemps eu une fascination par les motifs des pelages d'animaux (animal coat pattern en anglais). D'où provient la diversité des motifs des pelages d'animaux ? Pourquoi les léopards possèdent-ils des taches et ne possèdent pas des rayures ?

Toutes ces questions qui conduisent à comprendre le développement des motifs des pelages d'animaux chez les mammifères ont été un domaine d'études pour les biologistes et les chimistes. Une théorie sur la façon où ces motifs se produisent a été introduite par le mathématicien Alan Mathison Turing. Ce dernier a introduit un modèle pour lequel les motifs spatiaux peuvent se produire comme un résultat d'instabilité de la diffusion des substances chimiques morphogénétiques présentes dans la peau des animaux durant le stade embryonnaire du développement.



FIGURE 1.1 – Motifs du léopard, de girafe, du zèbre et du poisson.

1.1.1 Motivation biologique

La forme biologique d'un prémotif est inscrite dans les cellules : lors de la croissance de celles-ci, des produits chimiques (les morphogènes) se sont produites en concentration variable, et jouent suivant leur concentration le rôle d'activateur ou d'inhibiteur des pigments.

Pour l'obtention de tels motifs de pelages d'animaux, nous nous intéressons à un type spécifique de cellules appelées mélanocytes ; en effet, ce sont les cellules responsables de la production d'un pigment appelé la mélanine gouvernant la couleur de la peau.

Les biologistes pensent que les mélanocytes produisent de la mélanine basée sur la présence de certains activateurs et inhibiteurs des pigments. Chaque motif des pelages d'animaux est considéré comme étant le produit de certains activateurs et inhibiteurs [59].

1.1.2 Motivation mathématique (Principe de Turing)

En 1952, le mathématicien Alan Turing s'est intéressé à la morphogenèse et il a proposé un modèle basé sur des équations de **réaction–diffusion** (système de Turing) [76, 61]. Pour établir son modèle, Turing s'est appuyé sur la fameuse expérience réalisée par le biochimiste Belousov (1950) [62]. Ce dernier a utilisé des substances chimiques, et il a observé que la solution obtenue n'at-



FIGURE 1.2 – La réaction de Bolousov produisant des patterns temporels oscillants.

teint pas normalement l'état d'équilibre. En revanche il a observé qu'il y avait une apparition des patterns temporels oscillants jusqu'à épuisement d'un des réactifs.

De sa part, Turing a établi que, étant donné deux ensembles de substances chimiques, l'un est appelé **actif** et l'autre est appelé **inhibiteur** : si les deux ensembles interagissent et se diffusent dans un domaine fini et si en plus ces substances se propagent dans l'espace avec des différents taux de diffusions, alors nous pouvons générer des patterns spatiaux et stationnaires.

Pour bien comprendre l'idée séminale de Turing, voici l'image la plus frappante : supposons qu'un feu (l'actif) éclate dans une forêt sèche. Tout au début, il n'y a probablement pas de pompiers (les inhibiteurs) à proximité, mais avec leurs hélicoptères, ils peuvent dépasser le front de l'incen-

die et pulvériser sur les arbres des produits chimiques résistants au feu : quand le feu atteint les arbres traités, il s'éteint, l'incendie est arrêté.

Si plusieurs feux se déclarent spontanément et de façon aléatoire au sein de la forêt, plusieurs fronts d'incendie (des vagues d'activation) se propageront ; chaque incendie sera maîtrisé par les pompiers dans leurs hélicoptères (vagues d'inhibiteurs), mais il aura brûlé la forêt tout autour de son foyer initial. À la fin de la saison sèche, la forêt présentera des zones noires formées par les arbres brûlés et des zones vertes constituées d'arbres intacts.

Nous pouvons schématiser l'idée séminale de Turing comme suit : quand la diffusion n'est pas un paramètre essentiel c-à-d en absence de la diffusion (dans un milieu fortement brassé par exemple), les deux morphogènes réagissent et atteignent un état d'équilibre uniforme (pas de pattern). Dans le cas contraire, c-à-d en présence de la diffusion et sous certaines conditions appelées *conditions de Turing* (les morphogènes diffusent avec des vitesses différentes par exemple), une petite perturbation spatiale transforme l'état d'équilibre en un état instable engendrant des motifs spatiaux hétérogènes que nous appelons *Patterns de Turing*. En d'autres termes, sous les conditions de Turing la diffusion peut être déstabilisatrice et une petite perturbation spatiale peut être "instable" et engendrer un motif : une telle instabilité est appelée "diffusionnelle" ou bien "instabilité de turing".

1.1.3 Validation de l'idée de Turing

Après l'apparition de l'idée de Turing, plusieurs scientifiques ont travaillé pour vérifier cette idée, citons par exemple James Dickson Murray et Daniel Thomas. Murray s'est intéressé au modèle de Turing pour deux raisons. La première est que le chimiste Daniel Thomas (Université de Compiègne 1975) a travaillé sur le modèle de Turing et a pu démontrer par une expérience faite sur des produits chimiques spécifiques que le système de Turing peut engendrer des patterns spatiaux et stationnaires (patterns de Turing). Cette expérience était la première expérience confirmant que l'idée de Turing peut être fiable.

La deuxième raison pour laquelle Murray était intéressé par le modèle de Turing est que le modèle peut générer des patterns (des taches) comme ceux qui apparaissent sur la peau des animaux [61]. Murray a confirmé mathématiquement l'idée de Turing, et il a trouvé que, si en faisant juste varier le paramètre échelle de mesure γ (scale parameter) alors le modèle peut engendrer des patterns spatiaux qui prennent, suivant la valeur de γ , la forme de rayures ou de taches.

D'après ce qui précède, nous avons établi le rôle de la formation de patterns pour comprendre la formation de motifs de pelages d'animaux qui sont un résultat d'un principe de "réaction-diffusion" justifié par les chimistes : "activation, inhibition, diffusion". Un des exemples les plus connus et les plus populaires sur les systèmes de "réaction-diffusion" est la chimiotaxie.

1.1.4 La chimiotaxie

Un grand nombre des insectes et des animaux s'appuient sur une sensation aigüe de l'odorat, soit pour transmettre l'information entre les membres de l'espèce soit pour la prédation. Les produits chimiques impliqués dans ce processus sont appelés *phéromones*. Par exemple, la femelle du ver à soie libère une phéromone comme un attractif sexuel pour le mâle, qui dispose d'une antenne remarquablement efficace pour mesurer la concentration de la phéromone émise. Par conséquent, le mâle se déplace dans la direction de la forte et croissante concentration de la phéromone. La forte sensation de l'odorat pour la plupart des poissons est très importante pour la communication et la prédation ; par exemple, le requin dispose d'une sensation qui lui permet de sentir une goutte de sang diluée dans les tonnes d'eau, ce qui lui facilite la recherche de proie. L'exploitation la plus importante de la libération de phéromones est la direction du mouvement dirigé. Ce mouvement

dirigé par les produits chimiques est appelé la chimiotaxie, qui contrairement à la diffusion, dirige le mouvement vers un gradient de concentration.

Ce n'est pas seulement dans l'écologie des animaux et des insectes que la chimiotaxie est importante [43, 24, 68, 53, 21]. Elle peut être aussi cruciale dans les processus biologiques où il y a de nombreux exemples. Par exemple, quand une infection bactérienne envahit le corps, elle peut être attaquée par un mouvement de cellules vers la source par la suite de la chimiotaxie. Des preuves convaincantes suggèrent que les cellules de leucocytes dans le sang se dirigent vers la région d'inflammation bactérienne, pour l'attaquer, en se déplaçant suivant un gradient chimique causée par l'infection.

La modélisation théorique et mathématique de la chimiotaxie remonte aux travaux pionniers de Paltack [69] en 1950. Le modèle de Keller–Segel [51] introduit en 1970 reste le modèle le plus populaire et le plus connu pour le contrôle chimique des mouvement cellulaires. L'article [47] de Horstmann fournit une étude mathématique détaillée sur le modèle de Keller–Segel. Récemment, Bendahmane et al. [8] étudient l'existence des solutions pour un modèle de Keller–Segel dégénéré en tenant compte de l'effet de remplissage du volume.

L'analyse mathématique montre que les solutions du modèle de Keller–Segel peuvent modéliser des nombreux patterns spatiaux. Pour cela, l'étude de la formation de patterns pour le modèle de chimiotaxie est l'une des motivations principales de cette thèse.

Dans la suite, nous donnons un plan détaillé de la thèse avec les principaux résultats obtenus.

1.2 Plan de la thèse

Les travaux effectués dans cette thèse correspondent à l'étude théorique de la formation de patterns pour un modèle de chimiotaxie avec effet de remplissage du volume (volume-filling chemotaxis model en anglais) ; en plus, ils correspondent au développement d'une méthode numérique robuste et consistante capable de capturer la génération des patterns spatiaux pour le modèle étudié.

Dans le but d'obtenir un schéma numérique convergent sans restriction sur le maillage du domaine spatial ni sur les coefficients de transmissibilité, un schéma numérique est développé pour un modèle généralisé de Keller–Segel où nous prenons en compte une diffusion anisotrope et hétérogène de la densité cellulaire et de la concentration chimiotactique. Enfin, nous étudions l'existence de solutions faibles pour une équation parabolique dégénérée non linéaire découlant, soit de la modélisation de la chimiotaxie-fluide, soit de la modélisation d'un fluide compressible.

1.2.1 Chapitre 2 : Formation de patterns pour un modèle de chimiotaxie

Dans ce chapitre, nous nous intéressons à la formation de patterns spatiaux pour un système de réaction–convection–diffusion issue de la modélisation de l'effet de remplissage du volume pour un modèle de chimiotaxie. Ce système est donné par le couple des équations suivantes :

$$\begin{cases} \partial_t u - \operatorname{div} (d_1 (q(u) - q'(u)u) \nabla u) + \operatorname{div} (uq(u) \chi(v) \nabla v) = f(u, v), \\ \partial_t v - d_2 \Delta v = g(u, v). \end{cases} \quad (1.1)$$

Les deux équations du système (1.1) décrivent respectivement l'évolution en temps et en espace de la densité cellulaire u et de la concentration des substances chimiques v .

Le domaine occupé par les cellules et les substances chimiques est supposé fixe au cours du temps et est noté Ω , qui est un ensemble ouvert et borné de \mathbb{R}^2 . Nous posons $Q_{t_f} = \Omega \times (0, t_f)$, $\Sigma_{t_f} = \partial\Omega \times (0, t_f)$. Le système (1.1) est complété par la donnée de conditions initiales

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}) \geq 0, \quad v(\mathbf{x}, 0) = v_0(\mathbf{x}) \geq 0, \quad (1.2)$$

pour tout $\mathbf{x} \in \Omega$ et par la donnée de conditions aux limites de type Neumann homogène

$$(d_1 (q(u) - q'(u)u) \nabla u - uq(u) \chi(v) \nabla v) \cdot \mathbf{n} = 0, \quad \nabla v \cdot \mathbf{n} = 0, \quad (1.3)$$

pour tout $(\mathbf{x}, t) \in \Sigma_{t_f}$ où le vecteur \mathbf{n} désigne la normale à $\partial\Omega$ sortante de Ω . L'existence d'une solution globale pour le modèle de chimiotaxie (1.1) a été obtenue pour la première fois avec Painter et Hillen dans [44], où la prolifération des cellules n'était pas prise en compte. Plus tard, Wrzosek [79] a démontré l'existence d'une solution globale du système (1.1), pour n'importe quelle dimension de l'espace et avec une cinétique non nulle de la densité cellulaire f et une probabilité de pression q linéaire dans $[0, \bar{u}]$, où \bar{u} désigne le nombre total de cellules qui peuvent être logées dans n'importe quel site. Cette probabilité de pression reflète le fait que les cellules sont des particules solides ; elle est donnée par

$$q(u) = \begin{cases} 1 - \frac{u}{\bar{u}}, & 0 \leq u \leq \bar{u}, \\ 0, & u \geq \bar{u}. \end{cases} \quad (1.4)$$

Récemment, Wang et Hillen [78] ont montré l'existence d'une solution globale du système (1.1) pour une probabilité de pression plus réaliste (reflétant la propriété d'élasticité des cellules) en prenant une fonction non linéaire pour q .

Nous donnons les hypothèses portant sur le système (1.1) afin d'avoir des solutions globales ; elles sont données par

- (A1) d_1 et d_2 sont deux constantes positives.
- (A2) χ est une fonction de classe $C^2(\mathbb{R}, \mathbb{R})$ et telle que $\chi(v) > 0$.
- (A3) La probabilité de pression q est une fonction concave et décroissante de classe $C^3([0, \bar{u}])$ satisfaisant la condition suivante : il existe un nombre critique \bar{u} tel que $q(0) = 1$, $q(\bar{u}) = 0$, $0 < q(u) < 1$ pour $u \in (0, \bar{u})$ et $q(u) = 0$ pour tout $u > \bar{u}$.
- (A4) f est une fonction de classe $C^2(\mathbb{R} \times \mathbb{R})$ satisfaisant $f(0, v) \geq 0$ pour tout $v \geq 0$. En outre, il existe une constante $0 < u_c < \bar{u}$ telle que $f(u_c, v) = 0$ et $f(u, v) < 0$ pour tout $u > u_c$ et $v \geq 0$.
- (A5) g est une fonction de classe $C^2(\mathbb{R} \times \mathbb{R})$ satisfaisant $g(u, 0) \geq 0$ pour tout $u \geq 0$. En outre, il existe une constante $\bar{v} > 0$ telle que $g(u, \bar{v}) < 0$ pour tout $0 \leq u \leq \bar{u}$.

Analyse de stabilité linéaire pour le système (1.1)

Afin d'étudier la possibilité de générer des patterns spatiaux pour le système (1.1), nous appliquons l'idée séminale de Turing. Nous cherchons les états stationnaires stables (u_s, v_s) du système (1.1) en absence de toute variation spatiale ; ce qui amène à déterminer les états stationnaire stables pour le système suivant

$$\begin{aligned} \partial_t u &= f(u, v), & \partial_t v &= g(u, v) & \text{dans } Q_{t_f}, \\ f(u_s, v_s) &= g(u_s, v_s) = 0, & & & \text{dans } Q_{t_f}. \end{aligned} \quad (1.5)$$

Pour cette raison, nous effectuons une linéarisation autour de l'état stationnaire (u_s, v_s) (cette linéarisation est possible d'après les hypothèses de différentiabilité (A4)–(A5) des fonctions f et g) ; nous posons $\mathbf{w} = (u - u_s, v - v_s)$, ainsi l'équation (1.5) devient pour $|\mathbf{w}|$ assez petit,

$$\partial_t \mathbf{w} = A \mathbf{w}, \quad A = \begin{pmatrix} f_u & f_v \\ g_u & g_v \end{pmatrix}_{u_s, v_s},$$

où A est la matrice de stabilité (c-à-d la matrice jacobienne du système (1.5) calculée à l'état stationnaire (u_s, v_s)). Désormais, et sauf mention contraire, les dérivées partielles des fonctions f et g sont évaluées à l'état stationnaire (u_s, v_s) .

Par conséquent, les conditions nécessaires et suffisantes pour que l'état stationnaire homogène (u_s, v_s) soit linéairement stable, sont données par

$$\operatorname{tr}(A) = f_u + g_v < 0, \quad \det(A) = f_u g_v - f_v g_u > 0. \quad (1.6)$$

La deuxième étape de l'idée de Turing, consiste à déterminer les conditions pour lesquelles les états stationnaires stables (c-à-d ceux qui vérifient les conditions (1.6)) deviennent instables pour une petite perturbation spatiale. Pour cela, nous examinons une petite perturbation autour de l'état stationnaire (u_s, v_s) sous la forme suivante :

$$u = u_s + \varepsilon \tilde{u}(\mathbf{x}, t), \quad v = v_s + \varepsilon \tilde{v}(\mathbf{x}, t), \quad \varepsilon \ll 1. \quad (1.7)$$

En substituant la forme (1.7) dans le système (1.1), et en utilisant un développement asymptotique formel, nous obtenons le système suivant, après avoir assimilé les termes du premier ordre suivant ε et éliminé les tildes pour la simplicité

$$\begin{cases} \partial_t u = d_1 (q(u_s) - q'(u_s) u_s) \Delta u - u_s q(u_s) \chi(v_s) \Delta v + u f_u + v f_v, \\ \partial_t v = d_2 \Delta v + u g_u + v g_v. \end{cases} \quad (1.8)$$

Rappelons que les dérivées partielles des fonctions f et g figurant dans le système (1.8) sont des fonctions de u_s et v_s .

Une analyse sera étendue afin d'obtenir les conditions nécessaires concernant l'instabilité de l'état stationnaire (u_s, v_s) ; cette analyse s'appuie d'une part sur l'utilisation du principe de superposition des solutions dans le but de linéariser les termes spatiaux, et d'autre part, sur la détermination de la relation de dispersion associée au système (1.8) donné par l'équation $\lambda^2 + a(k^2) \lambda + b(k^2) = 0$, où $a(k^2) = (\xi d_1 + d_2) k^2 - (f_u + g_v)$ et $b(k^2) = \xi d_1 d_2 k^4 + (\psi g_u - d_2 f_u - \xi d_1 g_v) k^2 + f_u g_v - f_v g_u$.

ξ et ψ sont respectivement les coefficients de Δu et Δv de la première équation du système (1.8). Sous ces notations, les conditions nécessaires pour la générations de patterns spatiaux (conditions de Turing) pour le système (1.1) sont données par

$$\begin{aligned} f_u + g_v < 0, \quad f_u g_v - f_v g_u > 0, \quad \psi g_u - d_2 f_u - \xi d_1 g_v < 0, \\ (\psi g_u - d_2 f_u - \xi d_1 g_v)^2 - 4 \xi d_1 d_2 (f_u g_v - f_v g_u) > 0. \end{aligned} \quad (1.9)$$

Le système (1.1) dépend de plusieurs paramètres qui interviennent par exemple dans les fonctions cinétiques ou dans le terme de convection du système; ces paramètres sont appelés paramètres de bifurcation. La deuxième étape de l'étude de la formation de patterns, repose sur l'analyse de bifurcation pour localiser les éventuelles valeurs particulières des paramètres de bifurcations pour lesquels le comportement du système passe d'un état qualitatif à un autre. En effet, l'analyse de stabilité linéaire effectuée précédemment nous permet d'en déduire que la bifurcation se produit lorsque

$$\psi g_u - d_2 f_u - \xi d_1 g_v = -2 \sqrt{\xi d_1 d_2 (f_u g_v - f_v g_u)}. \quad (1.10)$$

Nous appliquons cette égalité pour déterminer les valeurs critiques de chaque paramètre de bifurcation. Par exemple, notons par $\vartheta = \chi(v_s)$ et remplaçons ξ par sa valeur dans l'équation (1.10) après avoir fixé tous les autres paramètres du système (1.1), nous obtenons la valeur critique ϑ_c donnée par

$$\vartheta_c = \frac{2 \sqrt{\xi d_1 d_2 (f_u g_v - f_v g_u)} - d_2 f_u - \xi d_1 g_v}{g_u u_s q(u_s)}.$$

D'après les conditions (1.9), la génération de patterns aura lieu lorsque $\vartheta > \vartheta_c$; ainsi, nous déterminons la gamme de nombres d'ondes instables associée au paramètre de bifurcation ϑ . Elle est définie par $k_1^2 < k^2 < k_2^2$, où k_1^2 et k_2^2 sont les zéros de l'équation $b(k^2) = 0$ données par

$$k_1^2 = \frac{\rho - \sqrt{\rho^2 - 4\xi d_1 d_2 (f_u g_v - f_v g_u)}}{2\xi d_1 d_2}, \quad k_2^2 = \frac{\rho + \sqrt{\rho^2 - 4\xi d_1 d_2 (f_u g_v - f_v g_u)}}{2\xi d_1 d_2},$$

avec $\rho = d_2 f_u + \xi d_1 g_v + \psi g_u$.

Dans le but de valider les résultats théoriques, la deuxième partie de ce chapitre est dédiée à l'examen numérique de la formation de patterns sous des paramètres spécifiques vérifiant les conditions de Turing (1.9). Nous utilisons la méthode de volumes finis pour discrétiser les équations du système (1.1)–(1.3) afin d'obtenir un schéma stable pour la discrétisation du terme de convection et vérifiant le principe de maximum (voir par exemple [78]) sur la solution (u, v) .

Schéma des volumes finis

Soit \mathcal{T} un maillage polygonal régulier et **admissible** (à préciser plus loin) du domaine Ω , constitué d'une famille de sous-domaines convexes, polygonaux non vides K de Ω appelés volumes de contrôle. Notons par h la taille de \mathcal{T} définie par $h = \max_{K \in \mathcal{T}} \text{diam}(K)$, où $\text{diam}(K)$ désigne la plus grande distance entre chaque paire de sommets du polygone convexe K .

Pour tout $K \in \mathcal{T}$, notons par \mathbf{x}_K le centre de K , $\mathcal{N}(K)$ l'ensemble de voisins de K . De plus, pour tout $L \in \mathcal{N}(K)$, nous notons par d_{KL} la distance entre \mathbf{x}_K et \mathbf{x}_L , par σ_{KL} l'interface entre K et son voisin L , et par η_{KL} la normale unitaire à σ_{KL} dirigée de K vers L . La figure 1.3 donne une illustration des volumes de contrôle ainsi que de la distance entre les centres.

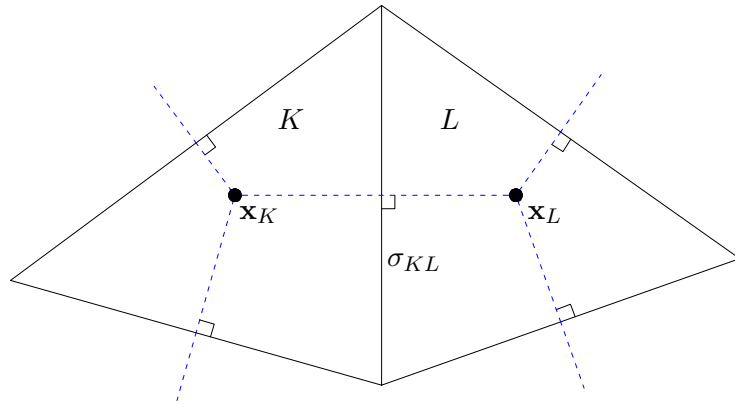


FIGURE 1.3 – Volumes de contrôle, centres, interfaces.

Pour tout $K \in \mathcal{T}$, nous désignons par $|K|$ la mesure de Lebesgue dans \mathbb{R}^2 du volume de contrôle K et par ∂K le bord de K formé par l'ensemble des arêtes de K .

\mathcal{T} est un maillage orthogonal admissible dans le sens de Eymard–Gallouët–Herbin [31].

L'admissibilité de \mathcal{T} implique que $\overline{\Omega} = \cup_{K \in \mathcal{T}} \overline{K}$, et pour tout volume de contrôle $L \in \mathcal{N}(K)$ (voisin du volume de contrôle K), le segment $[\mathbf{x}_K \mathbf{x}_L]$ joignant les centres de K et L est orthogonal à l'interface $\sigma_{KL} = \overline{K} \cap \overline{L}$.

Pour qu'un maillage conforme constitué de triangles vérifie la condition d'orthogonalité, il suffit de prendre le centre du volume de contrôle K , \mathbf{x}_K comme étant le centre du cercle circonscrit au triangle K . En outre, pour assurer que le centre \mathbf{x}_K soit dans le triangle K , nous imposons la condition que tous les angles du triangle K soient aigus.

L'avantage d'un tel maillage admissible, est de donner une approximation consistante de la dérivée dans la direction de la normale en utilisant uniquement deux points ; en effet par le développement de Taylor, nous déduisons que

$$\begin{aligned} \mathbf{x}_L - \mathbf{x}_K &= d_{KL} \eta_{KL}, \\ \nabla u(\mathbf{x}) \cdot \eta_{KL} &= \frac{u(\mathbf{x}_L) - u(\mathbf{x}_K)}{d_{KL}} + \mathcal{O}(h), \quad \forall \mathbf{x} \in \sigma = \overline{K} \cap \overline{L}. \end{aligned} \quad (1.11)$$

Une fois que le domaine est discrétisé, nous définissons l'espace discret $H_{\mathcal{T}}$ comme étant l'ensemble des fonctions constantes par morceaux sur les volumes de contrôle $K \in \mathcal{T}$. Chaque fonction $u_{\mathcal{T}} \in H_{\mathcal{T}}$ sera alors caractérisée par ses valeurs numériques $(u_K)_{K \in \mathcal{T}}$ telles que pour chaque volume de contrôle $K \in \mathcal{T}$, $u_{\mathcal{T}}|_K = u_K$. Plus précisément, la fonction discrète $u_{\mathcal{T}}$ s'écrit sous la forme d'une combinaison linéaire des fonctions caractéristiques $(\mathbf{1}_K)_{K \in \mathcal{T}}$ sous la forme suivante :

$$u_{\mathcal{T}}(\mathbf{x}) = \sum_{K \in \mathcal{T}} u_K \mathbf{1}_K(\mathbf{x}).$$

Nous donnons une définition du gradient discret $\nabla_{\mathcal{T}} u_{\mathcal{T}}$ qui est une fonction constante par diamant T_{KL} . Nous appelons un diamant T_{KL} associé à l'arête σ_{KL} , le polygone formé de quatre sommets \mathbf{x}_K , \mathbf{x}_L et les deux sommets de l'arête σ_{KL} (voir figure 1.4). Le gradient discret est défini par

$$\nabla_{\mathcal{T}} u_{\mathcal{T}}|_{T_{\sigma}} = \begin{cases} 2 \frac{u_L - u_K}{d_{KL}} \eta_{KL}, & \text{si } \sigma = \sigma_{KL}, \\ 0, & \text{si } \sigma \in \partial K \cap \partial \Omega. \end{cases}$$

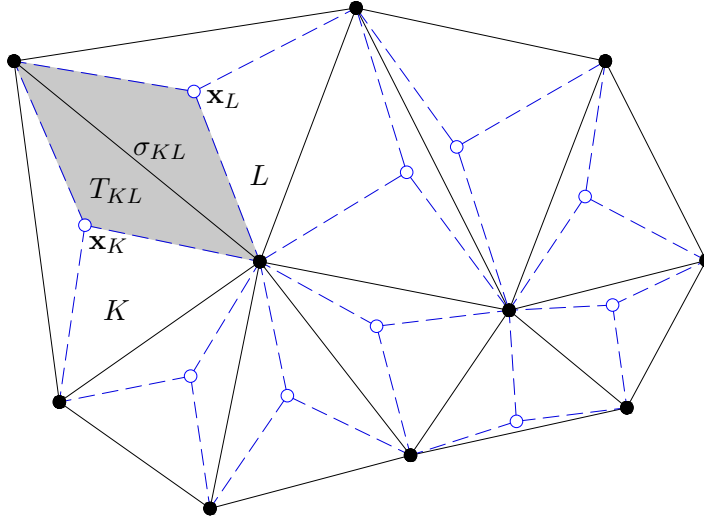


FIGURE 1.4 – Maillage admissible \mathcal{T} : volumes de contrôles, centres et diamants.

Ainsi, la norme dans $L^2(\Omega)$ de $\nabla_{\mathcal{T}} u_{\mathcal{T}}$ est donnée par

$$\|\nabla_{\mathcal{T}} u_{\mathcal{T}}\|_{L^2(\Omega)}^2 = \sum_{K \in \mathcal{T}} \sum_{L \in \mathcal{N}(K)} \frac{|\sigma_{KL}|}{d_{KL}} |u_K - u_L|^2.$$

La méthode des volumes finis consiste à intégrer directement les équations du modèle (contrairement à la méthode des éléments finis qui est basée sur les intégrales de la formulation variationnelle). Ainsi, les deux équations du système (1.1) sont intégrées sur chaque volume de contrôle $K \in \mathcal{T}$ et sur chaque intervalle de temps $]t_n, t_{n+1}[$, $n \in \mathbb{N}$ ($t_n = n\Delta t$, Δt est le pas de temps). Nous notons par A , Γ et Ψ les fonctions définies par $\Gamma(u) = q(u) - q'(u)u$, $A(u) = \int_0^u \Gamma(s) ds$ et $\Psi(u) = uq(u)\chi$.

Pour les équations du système (1.1), et en appliquant le théorème de la divergence cela donne

$$\begin{aligned} \int_{t_n}^{t_{n+1}} \int_K \partial_t u d\mathbf{x} dt - d_1 \sum_{L \in \mathcal{N}(K)} \int_{t_n}^{t_{n+1}} \int_{\sigma_{KL}} \nabla A(u) \cdot \eta_{KL} d\sigma(\mathbf{x}) dt \\ + \sum_{L \in \mathcal{N}(K)} \int_{t_n}^{t_{n+1}} \int_{\sigma_{KL}} \Psi(u) \nabla v \cdot \eta_{KL} d\sigma(\mathbf{x}) dt = \int_{t_n}^{t_{n+1}} \int_K f(u, v) d\mathbf{x} dt. \end{aligned} \quad (1.12)$$

et

$$\int_{t_n}^{t_{n+1}} \int_K \partial_t v d\mathbf{x} dt - d_2 \sum_{L \in \mathcal{N}(K)} \int_{t_n}^{t_{n+1}} \int_{\sigma_{KL}} \nabla v \cdot \eta_{KL} d\sigma(\mathbf{x}) dt = \int_{t_n}^{t_{n+1}} \int_K g(u, v) d\mathbf{x} dt, \quad (1.13)$$

où, $d\sigma(\mathbf{x})$ désigne la mesure de Lebesgue sur l'arête σ_{KL} .

Nous allons décrire brièvement une approximation de chaque terme des équations (1.12)–(1.13)

Les conditions initiales

$$u_K^0 = \frac{1}{|K|} \int_K u_0(\mathbf{x}) d\mathbf{x}, \quad v_K^0 = \frac{1}{|K|} \int_K v_0(\mathbf{x}) d\mathbf{x}.$$

Les termes d'évolution en temps

$$\frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} \int_K \partial_t w d\mathbf{x} dt \approx \frac{1}{\Delta t} \int_K (w(\mathbf{x}, t_{n+1}) - w(\mathbf{x}, t_n)) d\mathbf{x} \approx |K| \frac{w_K^{n+1} - w_K^n}{\Delta t}, \quad w \equiv u \text{ ou } v.$$

Les termes de diffusion Le maillage admissible a été considéré afin de donner une approximation simple des flux de diffusion ; en effet, l'équation (1.11) nous permet d'écrire

$$\frac{d_1}{\Delta t} \sum_{L \in \mathcal{N}(K)} \int_{t_n}^{t_{n+1}} \int_{\sigma_{KL}} \nabla A(u) \cdot \eta_{KL} d\sigma(\mathbf{x}) dt \approx d_1 \sum_{L \in \mathcal{N}(K)} \frac{|\sigma_{KL}|}{d_{KL}} (A(u_L^{n+1}) - A(u_K^{n+1})).$$

De la même manière, nous obtenons l'approximation du terme de diffusion de l'équation (1.13)

$$\frac{d_2}{\Delta t} \sum_{L \in \mathcal{N}(K)} \int_{t_n}^{t_{n+1}} \int_{\sigma_{KL}} \nabla v \cdot \eta_{KL} d\sigma(\mathbf{x}) dt \approx d_2 \sum_{L \in \mathcal{N}(K)} \frac{|\sigma_{KL}|}{d_{KL}} (v_L^{n+1} - v_K^{n+1}).$$

Le terme de convection L'approximation du terme de convection est différente de celle du terme de diffusion, ici nous voulons approcher le flux $\Psi(u) \nabla v \cdot \eta_{KL}$ (dépendant de u et de v) à l'interface σ_{KL} et à l'instant t_{n+1} . Le choix classique pour approcher un tel flux, consiste à utiliser un schéma décentré amont. Cette technique repose sur l'utilisation d'une fonction flux numérique G qui approche le flux $\Psi(u) \nabla v \cdot \eta_{KL}$ par le moyen des valeurs u_K , u_L et $dV_{KL} := \frac{|\sigma_{KL}|}{d_{KL}} (v_L - v_K)$. Nous donnons les propriétés données sur la fonction G d'arguments $(a, b, c) \in \mathbb{R}^3$:

- (i) Monotonie : $G(\cdot, b, c)$ est croissante pour tout $b, c \in \mathbb{R}$ et $G(a, \cdot, c)$ est décroissante pour tout $a, c \in \mathbb{R}$. Cette propriété est importante pour assurer le principe de maximum discret.
- (ii) Conservativité : $G(a, b, c) = -G(b, a, -c)$ pour tout $a, b, c \in \mathbb{R}$. Cette condition nous permet d'établir l'intégration par parties discrète essentielle pour la convergence du schéma numérique.
- (iii) Consistance : $G(a, a, c) = \Psi(a) c$ pour tout $a, c \in \mathbb{R}$.
- (iv) Majoration locale : Il existe une constante $C > 0$ telle que $|G(a, b, c)| \leq C(|a| + |b|)|c|$, pour tout $a, b, c \in \mathbb{R}$.
- (v) Continuité locale : Il existe un module de continuité $\omega : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ tel que $|G(a, b, c) - G(a', b', c)| \leq |c| \omega(|a - a'| + |b - b'|)$, pour tout $a, a', b, b', c \in \mathbb{R}$.

Nous prenons une fonction numérique G vérifiant les propriétés précédentes, par suite l'approximation du terme de convection est définie par

$$\frac{1}{\Delta t} \sum_{L \in \mathcal{N}(K)} \int_{t_n}^{t_{n+1}} \int_{\sigma_{KL}} \Psi(u) \nabla v \cdot \eta_{KL} d\sigma(\mathbf{x}) dt \approx \sum_{L \in \mathcal{N}(K)} G(u_K^{n+1}, u_L^{n+1}; dV_{KL}^{n+1})$$

En résumé, le schéma de type volumes finis proposé pour la discrétisation du problème (1.1)–(1.3) consiste à chercher $U = (u_K^{n+1})_{K \in \mathcal{T}, n \in [0, \dots, N]}$ et $V = (v_K^{n+1})_{K \in \mathcal{T}, n \in [0, \dots, N]}$ solution de

$$u_K^0 = \frac{1}{|K|} \int_K u_0(\mathbf{x}) d\mathbf{x}, \quad v_K^0 = \frac{1}{|K|} \int_K v_0(\mathbf{x}) d\mathbf{x}, \quad (1.14)$$

et

$$\begin{aligned} |K| \frac{u_K^{n+1} - u_K^n}{\Delta t} - d_1 \sum_{L \in \mathcal{N}(K)} \frac{|\sigma_{KL}|}{d_{KL}} (A(u_L^{n+1}) - A(u_K^{n+1})) \\ + \sum_{L \in \mathcal{N}(K)} G(u_K^{n+1}, u_L^{n+1}; dV_{KL}^{n+1}) = f(u_K^{n+1}, v_K^{n+1}), \quad (1.15) \\ |K| \frac{v_K^{n+1} - v_K^n}{\Delta t} - d_2 \sum_{L \in \mathcal{N}(K)} \frac{|\sigma_{KL}|}{d_{KL}} (v_L^{n+1} - v_K^{n+1}) = g(u_K^n, v_K^{n+1}). \end{aligned}$$

Sous les hypothèses (A1)–(A5) et supposons que $(u_0, v_0) \in (L^\infty(Q_{t_f}))^2$ et tel que $0 \leq u_0 \leq 1$, $v_0 \geq 0$. Alors le schéma (1.14)–(1.15) converge dans $L^2(Q_{t_f})$ quand $\Delta t, h \rightarrow 0$, pour une sous-suite, vers une solution faible (u, v) du problème (1.1)–(1.3).

La résolution de ce schéma nécessite toujours la résolution d'un système non linéaire. Pour cela, nous avons implémenté la méthode de Newton couplé avec l'algorithme du bigradient.

Enfin, des tests numériques en dimension deux sont faits pour simuler la formation de patterns pour le modèle (1.1)–(1.3) et pour montrer la robustesse du schéma (1.14)–(1.15) à capturer les patterns spatiaux.

1.2.2 Chapitre 3 : Un schéma volumes finis éléments finis pour capturer les patterns pour un modèle de chimiotaxie

Dans ce chapitre, nous nous sommes intéressés à un modèle de Keller–Segel similaire à celui de la section précédente, et pour lequel nous avons ajouté des tenseurs pour les termes de diffusion et

pour le terme de convection. Plus précisément, nous considérons le système suivant :

$$\begin{cases} \partial_t u - \operatorname{div} (\Lambda(\mathbf{x}) a(u) \nabla u - \Lambda(\mathbf{x}) \chi(u) \nabla v) = f(u) & \text{dans } Q_{t_f} = \Omega \times (0, t_f), \\ \partial_t v - \operatorname{div} (D(\mathbf{x}) \nabla v) = g(u, v) & \text{dans } Q_{t_f} = \Omega \times (0, t_f), \end{cases} \quad (1.16)$$

avec les conditions aux limites données sur le bord $\Sigma_T := \partial\Omega \times (0, t_f)$ par :

$$(\Lambda(\mathbf{x}) a(u) \nabla u - \Lambda(\mathbf{x}) \chi(u) \nabla v) \cdot \mathbf{n} = 0, \quad D(\mathbf{x}) \nabla v \cdot \mathbf{n} = 0, \quad (1.17)$$

où \mathbf{n} est le vecteur normal à $\partial\Omega$ sortant de Ω . Les conditions initiales sont données par :

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad v(\mathbf{x}, 0) = v_0(\mathbf{x}), \quad \mathbf{x} \in \Omega. \quad (1.18)$$

Le domaine occupé par les cellules et les substances chimiques est toujours noté par Ω , un ensemble ouvert et borné de \mathbb{R}^2 . Nous donnons également les hypothèses portant sur la diffusion de la densité cellulaire a , la sensibilité chimiotactique χ , les tenseurs de diffusion Λ et D , les cinétiques et les conditions initiales :

(A1) $a : [0, 1] \rightarrow \mathbb{R}$ est une fonction continue telle que : $a(0) \geq 0$, $a(1) \geq 0$ et $a(u) > 0$ pour tout $0 < u < 1$.

(A2) $\chi : [0, 1] \rightarrow \mathbb{R}$ est une fonction différentiable telle que : $\chi(0) = \chi(1) = 0$ et $\chi(u) > 0$ pour tout $0 < u < 1$.

(A3) Les tenseurs de diffusion Λ et D sont deux tenseurs symétriques bornés, uniformément positifs sur Ω , c-à-d tel que :

$$\forall \mathbf{w} \neq 0, \text{ il existe deux constantes strictement positives } T_- \text{ et } T_+ \text{ telles que } 0 < T_- |\mathbf{w}|^2 \leq \langle T(\mathbf{x}) \mathbf{w}, \mathbf{w} \rangle \leq T_+ |\mathbf{w}|^2 < \infty, T = \Lambda \text{ ou } D.$$

(A4) La fonction $f : \mathbb{R} \rightarrow \mathbb{R}$ est une fonction continue telle que : $f(0) \geq 0$ et $f(1) \leq 0$.

(A5) Les conditions initiales u_0 et v_0 sont deux fonctions dans $L^\infty(\Omega)$ tel que : $0 \leq u_0 \leq 1$ et $v_0 \geq 0$.

Notons par $A : \mathbb{R} \rightarrow \mathbb{R}$ la fonction Lipschitzienne et croissante définie par $A(u) = \int_0^u a(s) ds$, pour tout $u \in \mathbb{R}$.

Définition 1.1 (Solution faible). Sous les hypothèses (A1)–(A5). Le couple (u, v) est dit solution faible du système (1.16)–(1.17) si il vérifie :

$$\begin{aligned} 0 &\leq u(\mathbf{x}, t) \leq 1, \quad v(\mathbf{x}, t) \geq 0 \text{ pour presque tout } (\mathbf{x}, t) \in Q_{t_f}, \\ A(u) &\in L^2(0, t_f; H^1(\Omega)), \\ v &\in L^\infty(Q_{t_f}) \cap L^2(0, t_f; H^1(\Omega)), \end{aligned}$$

et pour tout $\varphi, \psi \in \mathcal{D}(\overline{\Omega} \times [0, t_f])$

$$\begin{aligned} & - \int_{\Omega} u_0(\mathbf{x}) \varphi(\mathbf{x}, 0) d\mathbf{x} - \iint_{Q_{t_f}} u \partial_t \varphi d\mathbf{x} dt + \iint_{Q_{t_f}} \Lambda(\mathbf{x}) \nabla A(u) \cdot \nabla \varphi d\mathbf{x} dt \\ & \quad - \iint_{Q_{t_f}} \Lambda(\mathbf{x}) \chi(u) \nabla v \cdot \nabla \varphi d\mathbf{x} dt = \iint_{Q_{t_f}} f(u) \varphi(\mathbf{x}, t) d\mathbf{x} dt, \\ & - \int_{\Omega} v_0(\mathbf{x}) \psi(\mathbf{x}, 0) d\mathbf{x} - \iint_{Q_{t_f}} v \partial_t \psi d\mathbf{x} dt \\ & \quad + \iint_{Q_{t_f}} D(\mathbf{x}) \nabla v \cdot \nabla \psi d\mathbf{x} dt = \iint_{Q_{t_f}} g(u, v) \psi d\mathbf{x} dt. \end{aligned}$$

L'intention de ce chapitre est de construire un schéma numérique robuste et convergent pour l'approximation du système (1.16)–(1.18). La difficulté repose sur la présence d'un tenseur anisotrope dans les termes de diffusion ; en effet, la méthode des volumes finis classique, malgré sa capacité à assurer la stabilité pour le terme de convection, ne permet pas la manipulation des problèmes de diffusion avec des tenseurs anisotropes, même si la condition d'orthogonalité est satisfaite. La raison de la non efficacité de la méthode de volumes finis, est que contrairement à l'approximation (1.11), il n'y a pas un moyen simple et direct pour approcher les flux de diffusion avec des tenseurs anisotropes et hétérogènes. Néanmoins, il est bien connu que la méthode des éléments finis nous permet de réaliser une discrétisation très simple des termes de diffusion avec des tenseurs pleins et sans imposer aucune restriction sur la géométrie du maillage ; par contre, des instabilités numériques peuvent se produire dans le cas d'une convection dominée. Une idée assez intuitive est donc de combiner les deux méthodes en donnant une discrétisation par la méthode des éléments finis conformes pour les termes de diffusion et une discrétisation par la méthode de volumes finis pour les autres termes. Par conséquent, nous construisons et nous étudions l'analyse de convergence d'un nouveau schéma appelé schéma CVFE afin de simuler et capturer la formation de patterns pour le modèle de chimiotaxie anisotrope.

Discrétisation du problème continu avec le schéma CVFE

Nous nous plaçons dans le cas de figure où les bordures du domaine Ω sont fixes au cours du temps. La construction de la solution approchée nécessite l'introduction de deux différentes discrétisations spatiales du domaine Ω , à savoir la *discrétisation primale* et la *discrétisation duale*. La discrétisation primale \mathcal{T} est une triangulation conforme du domaine Ω constituée d'un nombre fini de triangles qui forment une partition du domaine Ω c-à-d $\bar{\Omega} = \bigcup_{K \in \mathcal{T}} \bar{K}$ et $K \cap K' = \emptyset$ si $K \neq K'$. Pour chaque triangle $K \in \mathcal{T}$, nous notons par \mathbf{x}_K son centre de gravité, par h_K son diamètre, par ρ_K le diamètre du cercle circonscrit au triangle K et par $|K|$ la mesure de Lebesgue du triangle K . Ainsi, pour chaque sommet S du triangle K , nous notons par \mathcal{K}_S l'ensemble de triangles admettant S comme sommet et par \mathcal{E}_S l'ensemble des arêtes admettant S comme extrémité.

L'objectif de la technique CVFE est de reconstruire la solution approchée aux sommets du maillage primal. Pour ce faire, une nouvelle partition du domaine est définie de telle sorte que chaque sommet du maillage primal ne soit pas inclus que dans un triangle de la nouvelle partition. Le maillage issu de ce deuxième partitionnement est appelé *maillage dual* noté par \mathcal{M} et les polygones qui le composent *cellules duales* ou *volumes de contrôle duaux*. Pour chaque sommet S du triangle $K \in \mathcal{T}$, il existe un unique volume de contrôle dual M construit autour du sommet S en joignant les centres de gravité \mathbf{x}_K des triangles $K \in \mathcal{K}_S$ avec les milieux des arêtes $\sigma \in \mathcal{E}_S$. Le centre de M est noté par \mathbf{x}_M , il est confondu avec le sommet S autour duquel le volume de contrôle dual M est construit. Finalement, pour chaque paire de volumes de contrôle duaux M et M' , nous notons par $\sigma_{M,M'}^K$ le segment appartenant au triangle K et joignant le centre de gravité \mathbf{x}_K avec le milieu de l'arête $[\mathbf{x}_M \mathbf{x}_{M'}]$ où $K \in \mathcal{T}$ est un triangle tel que $M \cap K \neq \emptyset$ et $M' \cap K \neq \emptyset$. La figure 1.5 donne une illustration du maillage primal et du maillage dual correspondant.

Dans tout ce chapitre, nous limitons notre étude au cas d'une discrétisation uniforme en temps avec un pas de temps donné par $\Delta t = t_f / (N + 1)$, $N \in \mathbb{N}^*$ et par suite $t_n = n\Delta t$, $0 \leq n \leq N + 1$; ainsi, nous notons par $h = \max_{M \in \mathcal{M}} \text{diam}(M)$ le pas d'espace (la taille du maillage).

Pour chacune des deux reconstructions spatiales précédentes, nous définissons deux solutions approchées au sens du schéma CVFE :

- (i) Une solution volume fini $(u_{\mathcal{M},\Delta t}, v_{\mathcal{M},\Delta t})$ est définie comme étant une fonction constante

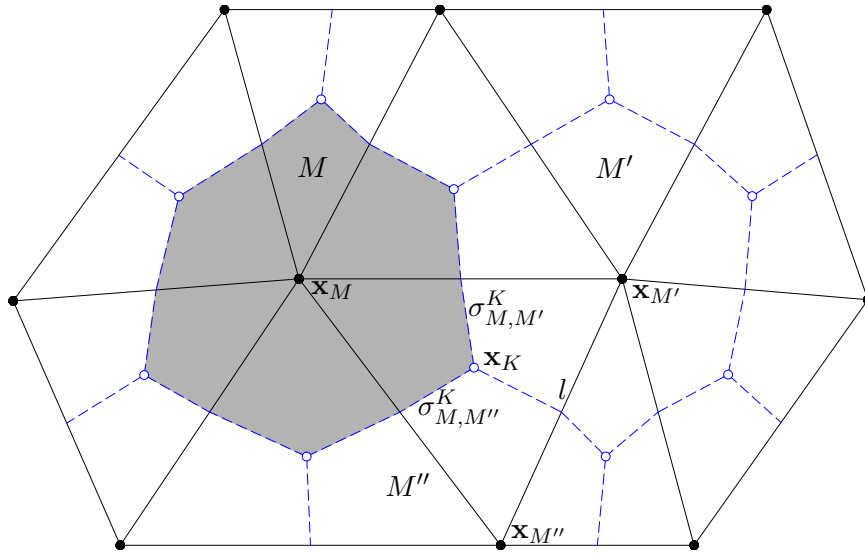


FIGURE 1.5 – Le maillage triangulaire primal \mathcal{T} et le maillage dual correspondant \mathcal{M} : volumes de contrôle, centres et interfaces.

par morceaux sur le maillage dual et telle que

$$\begin{aligned} (u_{\mathcal{M},\Delta t}(\mathbf{x}, 0), v_{\mathcal{M},\Delta t}(\mathbf{x}, 0)) &= (u_M^0, v_M^0) & \forall \mathbf{x} \in \overset{\circ}{M}, M \in \mathcal{M}, \\ (u_{\mathcal{M},\Delta t}(\mathbf{x}, t), v_{\mathcal{M},\Delta t}(\mathbf{x}, t)) &= (u_M^{n+1}, v_M^{n+1}) & \forall \mathbf{x} \in \overset{\circ}{M}, M \in \mathcal{M}, \forall t \in (t_n, t_{n+1}], \end{aligned}$$

où u_M^0 (resp. v_M^0) représente la valeur moyenne de la fonction u_0 (resp. v_0) sur M . L'espace discret de ces fonctions est noté $\mathcal{X}_{\mathcal{M},\Delta t}$.

- (ii) Une solution élément fini $v_{\mathcal{T},\Delta t}$ est définie comme étant une fonction continue et affine par triangle et telle que

$$\begin{aligned} v_{\mathcal{T},\Delta t}(\mathbf{x}, 0) &= v_{\mathcal{T}}^0(\mathbf{x}) & \forall \mathbf{x} \in \Omega, \\ v_{\mathcal{T},\Delta t}(\mathbf{x}, t) &= v_{\mathcal{T}}^{n+1}(\mathbf{x}) & \forall \mathbf{x} \in \Omega, \forall t \in (t_n, t_{n+1}], \end{aligned}$$

où $v_{\mathcal{T}}^{n+1}(\mathbf{x}) := \sum_{M \in \mathcal{M}} v_M^{n+1} \varphi_M(\mathbf{x})$ et $v_{\mathcal{T}}^0(\mathbf{x}) := \sum_{M \in \mathcal{M}} v_M^0 \varphi_M(\mathbf{x})$. L'espace discret de ces fonctions est noté $\mathcal{H}_{\mathcal{T},\Delta t}$. $(\varphi_M)_{M \in \mathcal{M}}$ est la base canonique de $\mathcal{H}_{\mathcal{T},\Delta t}$.

La fonction A est non linéaire, nous notons $A_{\mathcal{T},\Delta t} = A_{\mathcal{T}}(u_{\mathcal{T},\Delta t})$ la reconstruction éléments finis correspondante dans $\mathcal{H}_{\mathcal{T},\Delta t}$, et $A_{\mathcal{M},\Delta t} = A(u_{\mathcal{M},\Delta t})$ la reconstruction volumes finis correspondante dans $\mathcal{X}_{\mathcal{M},\Delta t}$. Pour les équations du système (1.16), et en appliquant la formule de Green–Gauss cela donne

$$\begin{aligned} \int_{t_n}^{t_{n+1}} \int_M \partial_t u \, d\mathbf{x} \, dt - \sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \Lambda \nabla A(u) \cdot \eta_{M,\sigma} \, d\sigma(\mathbf{x}) \, dt \\ + \sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \chi(u) \Lambda \nabla v \cdot \eta_{M,\sigma} \, d\sigma(\mathbf{x}) \, dt = \int_{t_n}^{t_{n+1}} \int_M f(u) \, d\mathbf{x} \, dt, \quad (1.19) \\ \int_{t_n}^{t_{n+1}} \int_M \partial_t v \, d\mathbf{x} \, dt - \sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} D \nabla v \cdot \eta_{M,\sigma} \, d\sigma(\mathbf{x}) \, dt = \int_{t_n}^{t_{n+1}} \int_M g(u, v) \, d\mathbf{x} \, dt, \end{aligned}$$

où, $\eta_{M,\sigma}$ est le vecteur normal unitaire à σ sortant de M et $d\sigma(\mathbf{x})$ désigne la mesure de Lebesgue sur l'arête σ .

Nous allons décrire brièvement l'approximation de chaque terme des équations du système (1.19).

Les conditions initiales

$$u_M^0 = \frac{1}{|M|} \int_M u_0(\mathbf{x}) d\mathbf{x}, \quad v_M^0 = \frac{1}{|M|} \int_M v_0(\mathbf{x}) d\mathbf{x}.$$

Les termes d'évolution en temps

$$\frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} \int_M \partial_t w d\mathbf{x} dt \approx \frac{1}{\Delta t} \int_M (w(\mathbf{x}, t_{n+1}) - w(\mathbf{x}, t_n)) d\mathbf{x} \approx |M| \frac{w_M^{n+1} - w_M^n}{\Delta t}, \quad w \equiv u \text{ ou } v.$$

Les termes de diffusion

$$\frac{1}{\Delta t} \sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \Lambda \nabla A(u) \cdot \eta_{M,\sigma} d\sigma(\mathbf{x}) dt \approx \sum_{\sigma \subset \partial M} \int_{\sigma} \Lambda \nabla A_{\mathcal{T}}(u_{\mathcal{T},\Delta t}(\mathbf{x}, t_{n+1})) \cdot \eta_{M,\sigma} d\sigma(\mathbf{x}).$$

Nous utilisons la reconstruction éléments finis $A_{\mathcal{T}}(u_{\mathcal{T},\Delta t})$ ainsi que la définition des fonctions de base $(\varphi_M)_{M \in \mathcal{M}}$ pour déduire l'approximation suivante

$$\sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \Lambda \nabla A(u) \cdot \eta_{M,\sigma} d\sigma(\mathbf{x}) dt \approx \Delta t \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} \Lambda_{M,M'}^K (A(u_{M'}^{n+1}) - A(u_M^{n+1})).$$

Le coefficient $\Lambda_{M,M'}^K$ est appelé le coefficient de transmissibilité défini par

$$\Lambda_{M,M'}^K = - \int_K \Lambda(\mathbf{x}) \nabla \varphi_M(\mathbf{x}) \cdot \nabla \varphi_{M'}(\mathbf{x}) d\mathbf{x}.$$

De la même manière, nous obtenons une approximation du terme de diffusion de l'équation en v

$$\sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \Lambda \nabla v \cdot \eta_{M,\sigma} d\sigma(\mathbf{x}) dt \approx \Delta t \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} D_{M,M'}^K (v_{M'}^{n+1} - v_M^{n+1}),$$

avec $D_{M,M'}^K = - \int_K D(\mathbf{x}) \nabla \varphi_M(\mathbf{x}) \cdot \nabla \varphi_{M'}(\mathbf{x}) d\mathbf{x}$.

Le terme de convection L'approximation du terme de convection est similaire à celle de la section précédente, nous voulons approcher le flux $\Lambda(\mathbf{x}) \chi(u_{\mathcal{M},\Delta t}) \nabla v_{\mathcal{T},\Delta t} \cdot \eta_{M,\sigma}$ à l'interface $\sigma \subset \partial M \cap \overline{K}$ et à l'instant t_{n+1} . Nous utilisons une fonction flux numérique G vérifiant les conditions (i)–(v) (données dans la section précédente) et qui approche le flux $\Lambda(\mathbf{x}) \chi(u_{\mathcal{M},\Delta t}) \nabla v_{\mathcal{T},\Delta t} \cdot \eta_{M,\sigma}$ par l'intermédiaire de u_M , $u_{M'}$ et l'approximation du gradient de v sur l'interface σ . En conclusion, l'approximation du terme de convection est donnée par

$$\Delta t \sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \chi(u) \Lambda \nabla v \cdot \eta_{M,\sigma} d\sigma(\mathbf{x}) dt \approx \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} |\sigma_{M,M'}^K| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K).$$

En résumé, le schéma CVFE proposé pour la discrétisation du problème (1.16)–(1.18) consiste à chercher $U = (u_M^{n+1})_{M \in \mathcal{M}, n \in [0, \dots, N]}$ et $V = (v_M^{n+1})_{M \in \mathcal{M}, n \in [0, \dots, N]}$ solution de

$$u_M^0 = \frac{1}{|M|} \int_M u_0(\mathbf{x}) \, d\mathbf{x}, \quad v_M^0 = \frac{1}{|M|} \int_M v_0(\mathbf{x}) \, d\mathbf{x}, \quad (1.20)$$

et

$$\begin{aligned} |M| \frac{u_M^{n+1} - u_M^n}{\Delta t} - \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} \Lambda_{M, M'}^K (A(u_{M'}^{n+1}) - A(u_M^{n+1})) \\ + \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} |\sigma_{M, M'}^K| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M, M'}^K) = |M| f(u_M^{n+1}), \quad (1.21) \\ |M| \frac{v_M^{n+1} - v_M^n}{\Delta t} - \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} D_{M, M'}^K (v_{M'}^{n+1} - v_M^{n+1}) = |M| g(u_M^n, v_M^{n+1}). \end{aligned}$$

Dans tout ce qui suit, nous considérons les hypothèses suivantes

(P1) $\Lambda_{M, M'}^K \geq 0$ et $D_{M, M'}^K \geq 0$, pour tout $M, M' \in \mathcal{M}$, et pour tout $K \in \mathcal{T}$.

(P2) il existe une constante $\kappa_{\mathcal{T}}$ strictement positive telle que : $\min_{K \in \mathcal{T}} \frac{|K|}{\text{diam}(K)^2} \geq \kappa_{\mathcal{T}}$.

Estimations discrètes, existence et convergence du schéma

Nous donnons des estimations a priori sur les solutions du schéma non linéaire (1.20)–(1.21) nécessaires pour établir l'existence d'une solution du schéma. Ensuite, nous présentons les résultats de compacité sur les solutions du schéma (translations en temps et en espace). Finalement, nous donnons le résultat principal de ce chapitre concernant la convergence du schéma basé sur l'utilisation du théorème de compacité de Kolmogorov.

Principe de maximum discret

Proposition 1.2. *Sous l'hypothèse (P1), supposons que $(u_M^n, v_M^n)_{M \in \mathcal{M}, n \in \{0, \dots, N+1\}}$ est une solution du schéma non linéaire (1.20)–(1.21). Alors, nous avons $0 \leq u_M^n \leq 1$ et $0 \leq v_M^n$ pour tout $M \in \mathcal{M}$, et pour tout $n \in \{0, \dots, N+1\}$. En outre, il existe une constante positive $\rho = \|v_0\|_{\infty} + \alpha t_f$, telle que $v_M^n \leq \rho$, pour tout $n \in \{0, \dots, N+1\}$.*

Nous donnons l'idée de la preuve pour u_M^{n+1} (la preuve pour v_M^{n+1} se traite de la même façon). Nous prenons un volume de contrôle dual M tel que $u_M^{n+1} = \min \{u_{M'}^{n+1}\}_{M' \in \mathcal{M}}$ et multiplions la première équation du système (1.21) par $-(u_M^{n+1})^-$ et faisons la sommation sur $M \in \mathcal{M}$. Nous obtenons la positivité de u_M^{n+1} en utilisant la formule d'intégration par parties discrètes, l'hypothèse sur la positivité des coefficients de transmissibilité et par le prolongement continue par zéro des fonctions $\chi(u)$ et $f(u)$ pour $u \leq 0$. D'une manière similaire, et en multipliant la première équation du système (1.21) par $(u_M^{n+1} - 1)^+$, nous obtenons que $u_M^{n+1} \leq 1$.

Estimations a priori Maintenant, nous donnons des propriétés discrètes sur le schéma non linéaire (1.20)–(1.21).

Proposition 1.3. *Sous les hypothèses (P1) et (P2), supposons que $(u_M^n, v_M^n)_{M \in \mathcal{M}, n \in \{0, \dots, N+1\}}$ est une solution du schéma non linéaire (1.20)–(1.21). Alors, il existe une constante $C > 0$ (indépendante de h et Δt) et telle que*

$$\sum_{M \in \mathcal{M}} |M| (v_M^{n+1})^2 + \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} D_{M,M'}^K (v_M^{n+1} - v_{M'}^{n+1})^2 \leq \sum_{M \in \mathcal{M}} |M| (v_M^0)^2 \leq C,$$

pour tout $n \in \{0, \dots, N+1\}$.

Notons par $\mathcal{B} : \mathbb{R} \rightarrow \mathbb{R}$ la fonction définie par $\mathcal{B}(u) = \int_0^u A(s) ds$. Nous avons une estimation similaire à celle de la Proposition 1.3, elle est donnée par la proposition suivante

Proposition 1.4. *Pour tout $n \in \{0, \dots, N+1\}$, il existe une constante $C > 0$ (indépendante de h et Δt) et telle que*

$$\sum_{M \in \mathcal{M}} |M| \mathcal{B}(u_M^{n+1}) + \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (A(u_M^{n+1}) - A(u_{M'}^{n+1}))^2 \leq C,$$

En utilisant les propositions 1.3 et 1.4, le théorème du point fixe de Brouwer assure l'existence d'au moins une solution $(u_M^n, v_M^n)_{M \in \mathcal{M}, n \in \{0, \dots, N+1\}}$ du problème discret non linéaire (1.20)–(1.21).

Estimations de compacité sur la solution discrète Nous donnons ici les translations en temps et en espace nécessaires pour montrer la convergence du schéma (1.20)–(1.21). En particulier, nous avons la proposition suivante :

Proposition 1.5. *Sous les hypothèses (P1) et (P2), il existe une constante $C > 0$ indépendante de h et τ et telle que*

$$\iint_{\Omega \times (0, t_f - \tau)} |w_{\mathcal{M}_h, \Delta t}(t + \tau, \mathbf{x}) - w_{\mathcal{M}_h, \Delta t}(\mathbf{x}, t)|^2 d\mathbf{x} dt \leq C(\tau + \Delta t), \quad \text{pour tout } \tau \in (0, t_f),$$

et

$$\int_0^{t_f} \int_{\Omega'} |w_{\mathcal{M}_h, \Delta t}(\mathbf{x} + \mathbf{y}, t) - w_{\mathcal{M}_h, \Delta t}(\mathbf{x}, t)| d\mathbf{x} dt \leq C(|\mathbf{y}| + h), \quad \text{pour tout } \mathbf{y} \in \mathbb{R}^2,$$

avec $\Omega' = \{\mathbf{x} \in \Omega, [\mathbf{x}, \mathbf{x} + \mathbf{y}] \subset \Omega\}$ et $w_{\mathcal{M}_h, \Delta t} = A(u_{\mathcal{M}_h, \Delta t})$ ou $v_{\mathcal{M}_h, \Delta t}$.

Avant de passer au résultat de convergence, nous montrons un lemme essentiel qui assure que la reconstruction éléments finis et la reconstruction volumes finis se rapprochent l'une de l'autre quand la taille du maillage h tend vers zéro. Ce résultat est donné par le lemme suivant :

Lemma 1.6. *Les suites $(A(u_{\mathcal{M}_h, \Delta t}) - A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t}))_{h, \Delta t}$ et $(v_{\mathcal{M}_h, \Delta t} - v_{\mathcal{T}_h, \Delta t})_{h, \Delta t}$ convergent fortement vers zéro dans $L^2(Q_T)$ quand $h \rightarrow 0$.*

Nous sommes en mesure d'annoncer le résultat principal de ce chapitre concernant la convergence du schéma non linéaire (1.20)–(1.21) vers une solution faible du problème continue (1.16)–(1.18) et quand les pas de discrétisations spatio-temporel tendent vers zéro.

Theorem 1.7. *Soit $(\mathcal{T}_m)_{m \geq 1}$ une suite de triangulations conformes de Ω telles que $h_{\mathcal{T}_m} \rightarrow 0$ quand $m \rightarrow \infty$, et soit $(\Delta t_m)_{m \geq 1}$ une suite de pas de temps telle que $\Delta t_m \rightarrow 0$ quand $m \rightarrow \infty$. Sous les hypothèses (P1) et (P2), la solution discrète $(u_{\mathcal{M}_m, \Delta t_m}, v_{\mathcal{M}_m, \Delta t_m})_m$ converge fortement dans $L^q(Q_{t_f})$ pour tout $q \in [1, \infty)$ vers une solution faible du problème continue (1.16)–(1.18) au sens de la définition (1.1) quand $m \rightarrow \infty$.*

La résolution du schéma (1.20)–(1.21) nécessite toujours la résolution d'un système non linéaire. Pour cela, la méthode de Newton est employée pour approcher la solution U^{n+1} de l'équation non linéaire définie par (1.20). Cet algorithme est couplé par une méthode de bigradient conjugué pour résoudre les systèmes linéaires découlant de la méthode de Newton aussi bien que la solution V^{n+1} de l'équation linéaire définie par (1.21).

La dernière partie de ce chapitre est constituée de trois simulations numériques effectuées avec le schéma numérique (1.20)–(1.21). Ces simulations ont pour but de vérifier l'efficacité du schéma proposé pour capturer la génération des patterns spatiaux pour le modèle de chimiotaxie sous l'effet du remplissage du volume. La première simulation consiste à prendre des tenseurs isotropes (c-à-d proportionnels à la matrice d'identité) ; ce cas d'étude est important afin de vérifier que les résultats effectués avec le schéma (1.20)–(1.21) sont cohérents avec ceux effectués en utilisant le schéma de volumes finis classique. La deuxième simulation consiste à prendre Λ comme étant un tenseur anisotrope diagonale dont les éléments diagonaux sont distincts, et à prendre une matrice isotrope pour D . Ce test montre que le modèle (1.16)–(1.18) génère des patterns spatiaux qui se propagent dans la direction d'un axe faisant avec l'axe des abscisses un angle identique à celui calculé en utilisant la matrice de passage pour le tenseur Λ . La dernière simulation consiste aussi à prendre une matrice isotrope pour D et à prendre une matrice anisotrope hétérogène pour Λ . Pour ce dernier test, nous obtenons également des patterns spatiaux pour lesquels la direction de propagation est cohérente avec l'angle calculé par la matrice de passage (matrice de rotation) du tenseur Λ .

1.2.3 Chapitre 4 : Un schéma CVFE non linéaire pour le modèle de Keller–Segel modifié

Dans ce chapitre, nous nous sommes intéressés à un modèle généralisé de Keller–Segel anisotrope et dégénéré. Le modèle ressemble à celui de la section précédente, pour lequel nous avons ajouté des tenseurs pour les termes de diffusion et pour le terme de convection. Plus précisément, nous considérons la formulation suivante :

$$\begin{cases} \partial_t u - \operatorname{div} (\Lambda(\mathbf{x}) a(u) \nabla u - \Lambda(\mathbf{x}) \chi(u) \nabla v) = f(u) & \text{dans } Q_{t_f} = \Omega \times (0, t_f), \\ \partial_t v - \operatorname{div} (D(\mathbf{x}) \nabla v) = g(u, v) & \text{dans } Q_{t_f} = \Omega \times (0, t_f), \end{cases} \quad (1.22)$$

avec les conditions aux limites données sur le bord $\Sigma_{t_f} := \partial\Omega \times (0, t_f)$ par

$$(\Lambda(\mathbf{x}) a(u) \nabla u - \Lambda(\mathbf{x}) \chi(u) \nabla v) \cdot \mathbf{n} = 0, \quad D(\mathbf{x}) \nabla v \cdot \mathbf{n} = 0, \quad (1.23)$$

et les conditions initiales données par :

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad v(\mathbf{x}, 0) = v_0(\mathbf{x}), \quad \mathbf{x} \in \Omega. \quad (1.24)$$

Le domaine occupé par les cellules et les substances chimiques est toujours noté Ω , un ensemble ouvert et borné de \mathbb{R}^2 . Nous donnons également les hypothèses portant sur la diffusion de la densité cellulaire dégénérée a , la sensibilité chimiotactique χ , les tenseurs de diffusion Λ et D , les cinétiques et les conditions initiales.

(A1) $a : [0, 1] \rightarrow \mathbb{R}$ est une fonction continue telle que $a(0) = a(1) = 0$ et $a(u) > 0$ pour tout $0 < u < 1$.

(A2) $\chi : [0, 1] \rightarrow \mathbb{R}$ est une fonction dérivable telle que $\chi(0) = \chi(1) = 0$ et $\chi(u) > 0$ pour tout $0 < u < 1$. En outre, il existe une fonction $\mu \in \mathcal{C}^0([0, 1], \mathbb{R})$ telle que $\mu(u) = \frac{\chi(u)}{a(u)}$, $\mu(0) = \mu(1) = 0$ et $\mu(u) > 0$ pour tout $0 < u < 1$.

(A3) Les tenseurs de diffusion Λ et D sont deux tenseurs symétriques bornés, uniformément positifs sur Ω , c-à-d tels que :

$$\forall \mathbf{w} \neq 0, \text{ il existe deux constantes strictement positives } T_- \text{ et } T_+ \text{ telles que } 0 < T_- |\mathbf{w}|^2 \leq \langle T(\mathbf{x})\mathbf{w}, \mathbf{w} \rangle \leq T_+ |\mathbf{w}|^2 < \infty, T = \Lambda \text{ ou } D.$$

(A4) La fonction $f : \mathbb{R} \rightarrow \mathbb{R}$ est une fonction continue telle que $f(0) \geq 0$ et $f(1) \leq 0$.

(A5) Les conditions initiales u_0 et v_0 sont deux fonctions dans $L^\infty(\Omega)$ telles que $0 \leq u_0 \leq 1$ et $v_0 \geq 0$.

Notons $\xi : [0, 1] \rightarrow \mathbb{R}$ la fonction Lipschitzienne et croissante définie par $\xi(u) = \int_0^u \sqrt{a(s)} ds$, pour tout $u \in \mathbb{R}$.

Definition 1.8 (Solution faible). Sous les hypothèses (A1)–(A5), le couple (u, v) est dit solution faible du système (1.22)–(1.24) si il vérifie :

$$\begin{aligned} 0 &\leq u(\mathbf{x}, t) \leq 1, \quad v(\mathbf{x}, t) \geq 0 \text{ pour presque tout } (\mathbf{x}, t) \in Q_{t_f}, \\ \xi(u) &\in L^2(0, t_f; H^1(\Omega)), \\ v &\in L^\infty(Q_{t_f}) \cap L^2(0, t_f; H^1(\Omega)), \end{aligned}$$

et pour tout $\varphi, \psi \in \mathcal{D}(\overline{\Omega} \times [0, t_f])$

$$\begin{aligned} & - \int_{\Omega} u_0(\mathbf{x}) \varphi(\mathbf{x}, 0) d\mathbf{x} - \iint_{Q_{t_f}} u \partial_t \varphi d\mathbf{x} dt + \iint_{Q_{t_f}} \sqrt{a(u)} \Lambda(\mathbf{x}) \nabla \xi(u) \cdot \nabla \varphi d\mathbf{x} dt \\ & \quad - \iint_{Q_{t_f}} \Lambda(\mathbf{x}) \chi(u) \nabla v \cdot \nabla \varphi d\mathbf{x} dt = \iint_{Q_{t_f}} f(u) \varphi(\mathbf{x}, t) d\mathbf{x} dt, \\ & - \int_{\Omega} v_0(\mathbf{x}) \psi(\mathbf{x}, 0) d\mathbf{x} - \iint_{Q_{t_f}} v \partial_t \psi d\mathbf{x} dt \\ & \quad + \iint_{Q_{t_f}} D(\mathbf{x}) \nabla v \cdot \nabla \psi d\mathbf{x} dt = \iint_{Q_{t_f}} g(u, v) \psi(\mathbf{x}, t) d\mathbf{x} dt. \end{aligned}$$

Dans ce chapitre, nous construisons et analysons un schéma numérique pour la discrétisation du système (1.22)–(1.24). Nous adoptons le principe du schéma non linéaire CVFE proposé dans la section précédente pour la construction du schéma. Dans ce dernier schéma, les degrés de liberté sont affectés aux sommets d'un maillage triangulaire primal, comme dans la méthode des éléments finis, tandis que les équations de bilan sont discrétisées sur un maillage dual spécifique (le maillage barycentrique dual appelé encore le maillage de Donald) en utilisant les flux de diffusion fournis par la reconstruction des éléments finis sur le maillage triangulaire primal. Plus particulièrement, les termes de diffusion sont approchés par l'intermédiaire des flux de diffusion découlant de la reconstruction des éléments finis et en utilisant un schéma de Godunov. Par contre, le terme de convection est approché par le moyen d'un schéma amont original ; en effet, la fonction χ est définie comme étant le produit de deux fonctions μ et a , le schéma proposé consiste à prendre le schéma amont classique d'une part, pour la discrétisation de la fonction μ qui sera approchée par l'intermédiaire d'une fonction flux numérique et à prendre le schéma de Godunov d'autre part, pour approcher la fonction a qui intervient dans la fonction de la sensibilité chimiotactique χ . Les approximations utilisées dans ce schéma sont cruciales afin d'assurer le principe de maximum, la stabilité et la convergence du schéma sans aucune restriction sur les coefficients de transmissibilité.

Comme dans la section précédente, la discrétisation du système (1.22)–(1.24) exige la construction de deux types d'approximations : l'approximation éléments finis sur un maillage triangulaire primal et une approximation volumes finis sur le maillage barycentrique correspondant.

Nous construisons le maillage primal \mathcal{T} à partir d'un nombre fini de triangles disjoints qui forment une partition du domaine Ω supposé polygonal ; nous avons alors $\bigcup_{T \in \mathcal{T}} \bar{T} = \bar{\Omega}$, et $T \cap T' = \emptyset$ si $T \neq T'$. Nous notons par \mathcal{V} l'ensemble de tous les sommets de la triangulation conforme \mathcal{T} et par \mathcal{E} l'ensemble de toutes les arêtes de la triangulation \mathcal{T} . Pour chaque sommet $K \in \mathcal{V}$ (localisé à la position \mathbf{x}_K), nous notons par \mathcal{E}_K le sous-ensemble de \mathcal{E} constitué des arêtes admettant le sommet K comme une extrémité. Une arête joignant deux sommets K et L est notée par σ_{KL} .

Pour la construction du maillage barycentrique dual, nous notons par \mathcal{T}_K l'ensemble de tous les triangles admettant K comme un sommet. Il existe un volume de contrôle unique ω_K associé au sommet K , il est construit autour de ce sommet en joignant les centres de gravité \mathbf{x}_T des triangles $T \in \mathcal{T}_K$ avec les milieux \mathbf{x}_σ des arêtes $\sigma \in \mathcal{E}_K$ (voir Figure 1.6). Nous déterminons les constantes

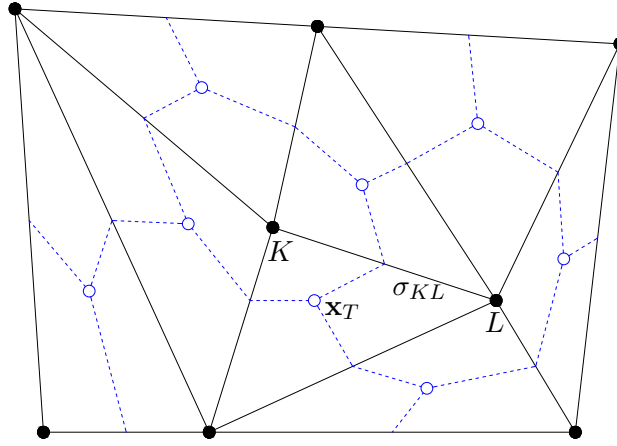


FIGURE 1.6 – Le maillage triangulaire primal \mathcal{T} et le maillage de Donald dual \mathcal{M} .

$h_{\mathcal{T}}$ et $\theta_{\mathcal{T}}$ comme étant la taille et la régularité du maillage Ω respectivement. Elles sont définies par :

$$h_{\mathcal{T}} = \max_{T \in \mathcal{T}}, \quad \theta_{\mathcal{T}} = \max_{T \in \mathcal{T}} \frac{h_T}{\rho_T},$$

où h_T est le diamètre du triangle T et ρ_T est le diamètre du cercle circonscrit au triangle T .

Nous notons $\mathcal{H}_{\mathcal{T}}$ l'espace usuel des éléments finis c-à-d l'espace \mathbb{P}_1 -discret formé des fonctions continues et affines par morceaux sur le maillage triangulaire primal

$$\mathcal{H}_{\mathcal{T}} = \{\phi \in C^0(\bar{\Omega}) ; \phi|_T \in \mathbb{P}_1(\mathbb{R}), \forall T \in \mathcal{T}\} \subset H^1(\Omega),$$

et par $(\varphi_K)_{K \in \mathcal{V}}$ sa base canonique définie par $\varphi_K(\mathbf{x}_L) = \delta_{KL}$ (δ est le symbole de Kronecker). En outre, nous considérons l'espace discret de volumes finis $\mathcal{X}_{\mathcal{M}}$ formé des fonctions constantes par morceaux sur le maillage barycentrique dual

$$\mathcal{X}_{\mathcal{M}} = \{\phi : \Omega \longrightarrow \bar{\mathbb{R}}, \phi|_{\omega_K} \text{ est une constante}, \forall K \in \mathcal{V}\}.$$

Dans le but de simplifier les notations dans le chapitre, nous limitons notre étude au cas d'une discrétisation uniforme en temps. Pour cela, nous notons $\Delta t = t_f/(N+1)$ le pas de temps uniforme et $t_n = n\Delta t$ de sorte que $t^0 = 0$ et $t_{N+1} = t_f$. Une fois l'intervalle de temps discrétisé, nous introduisons les espaces discrets spatiaux et temporels définis par :

$$\begin{aligned} \mathcal{H}_{\mathcal{T}, \Delta t} &= \{\phi \in L^2(0, t_f; H^1(\Omega)) ; \phi(\cdot, t) = \phi(\cdot, t_{n+1}) \in \mathcal{H}_{\mathcal{T}}, \forall t \in (t_n, t_{n+1}], 0 \leq n \leq N\}, \\ \mathcal{X}_{\mathcal{M}, \Delta t} &= \{\phi \in L^\infty(Q_{t_f}) ; \phi(\cdot, t) = \phi(\cdot, t_{n+1}) \in \mathcal{X}_{\mathcal{M}}, \forall t \in (t_n, t_{n+1}], 0 \leq n \leq N\}. \end{aligned}$$

Pour un vecteur donné $(u_K^n)_{n \in \{0, \dots, N+1\}, K \in \mathcal{V}}$ (resp. $(v_K^n)_{n \in \{0, \dots, N+1\}, K \in \mathcal{V}}$), il existe une unique fonction $u_{\mathcal{T}, \Delta t} \in \mathcal{H}_{\mathcal{T}, \Delta t}$ (resp. $v_{\mathcal{T}, \Delta t} \in \mathcal{H}_{\mathcal{T}, \Delta t}$) et une unique fonction $u_{\mathcal{M}, \Delta t} \in \mathcal{X}_{\mathcal{M}, \Delta t}$ (resp. $v_{\mathcal{M}, \Delta t} \in \mathcal{X}_{\mathcal{M}, \Delta t}$) telles que :

$$\begin{aligned} u_{\mathcal{T}, \Delta t}(\mathbf{x}_K, t_{n+1}) &= u_{\mathcal{M}, \Delta t}(\mathbf{x}_K, t_{n+1}) = u_K^{n+1}, & \forall K \in \mathcal{V}, \forall n \in \{0, \dots, N\}, \\ v_{\mathcal{T}, \Delta t}(\mathbf{x}_K, t_{n+1}) &= v_{\mathcal{M}, \Delta t}(\mathbf{x}_K, t_{n+1}) = v_K^{n+1}, & \forall K \in \mathcal{V}, \forall n \in \{0, \dots, N\}. \end{aligned}$$

Soit m_K la mesure de Lebesgue dans \mathbb{R}^2 du volume de contrôle ω_K . Pour chaque couple $(K, L) \in \mathcal{V}^2$, nous notons Λ_{KL} et D_{KL} les coefficients de transmissibilité définis par :

$$\mathbf{T}_{KL} = \int_{\Omega} \mathbf{T}(\mathbf{x}) \nabla \varphi_K(\mathbf{x}) \cdot \nabla \varphi_L(\mathbf{x}) d\mathbf{x}, \quad \mathbf{T} \equiv \Lambda \text{ ou } D. \quad (1.25)$$

Le schéma proposé consiste à intégrer le système (1.22) sur $\omega_K \times [t_n, t_{n+1}]$ avec $K \in \mathcal{V}$ et à utiliser le théorème de Green–Gauss pour passer à l’interface de chaque volume dual.

Discretisation de la deuxième équation du système (1.22) Nous détaillons la discrétisation de la première équation du système (1.22). Le théorème de Green–Gauss donne :

$$\begin{aligned} \int_{t_n}^{t_{n+1}} \int_{\omega_K} \partial_t u d\mathbf{x} dt - \sum_{\sigma \subset \partial \omega_K} \int_{t_n}^{t_{n+1}} \int_{\sigma} \Lambda(\mathbf{x}) a(u) \nabla u \cdot \eta_{\sigma} d\sigma(\mathbf{x}) dt \\ + \sum_{\sigma \subset \partial \omega_K} \int_{t_n}^{t_{n+1}} \int_{\sigma} \mu(u) \Lambda(\mathbf{x}) a(u) \nabla v \cdot \eta_{\sigma} d\sigma(\mathbf{x}) dt = \int_{t_n}^{t_{n+1}} \int_{\omega_K} f(u) d\mathbf{x} dt, \end{aligned}$$

Le terme d’évolution en temps

$$\frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} \int_{\omega_K} \partial_t u d\mathbf{x} dt \approx \frac{1}{\Delta t} \int_{\omega_K} (u(\mathbf{x}, t_{n+1}) - u(\mathbf{x}, t_n)) d\mathbf{x} \approx m_K \frac{u_K^{n+1} - u_K^n}{\Delta t}.$$

Le terme de diffusion Nous voulons approcher le flux $\mathcal{F}(u) = -\Lambda(\mathbf{x}) a(u) \nabla u \cdot \eta_{\sigma}$ sur l’interface σ (ici σ est l’interface partagée par ω_K avec son volume de contrôle dual voisin ω_L et appartenant au triangle T tel que $T \cap \omega_K \neq \emptyset \neq T \cap \omega_L$). Le schéma de Godunov consiste à approcher le flux \mathcal{F} en utilisant une fonction flux numérique \mathcal{F}_{KL} définie par :

$$\mathcal{F}_{KL}(u_K, u_L) = \begin{cases} \min_{u \in [u_K, u_L]} \Lambda_{KL}(u_K - u_L) a(u), & \text{si } u_K \leq u_L, \\ \max_{u \in [u_L, u_K]} \Lambda_{KL}(u_K - u_L) a(u), & \text{si } u_L \leq u_K, \end{cases}$$

où $\Lambda_{KL}(u_K - u_L)$ représente une approximation du flux de diffusion $-\Lambda \nabla u \cdot \eta_{\sigma}$ sur l’interface σ et provenant de la construction éléments finis (comme pour la section précédente) et où $\Lambda_{KL} = \sum_{K \in \mathcal{T}} \Lambda_{KL}^T$ avec $\Lambda_{KL}^T = \int_T \Lambda \nabla \varphi_K \cdot \nabla \varphi_L d\mathbf{x}$.

En utilisant les propriétés classiques sur les fonctions Lipschitziennes min et max, nous obtenons

$$\mathcal{F}_{KL}(u_K, u_L) = \begin{cases} \Lambda_{KL}(u_K - u_L) \max_{u \in [u_K, u_L]} a(u), & \text{si } \Lambda_{KL} \geq 0, \\ \Lambda_{KL}(u_K - u_L) \min_{u \in [u_K, u_L]} a(u), & \text{si } \Lambda_{KL} \leq 0, \\ \Lambda_{KL}(u_K - u_L) \max_{u \in [u_L, u_K]} a(u), & \text{si } \Lambda_{KL} \geq 0, \\ \Lambda_{KL}(u_K - u_L) \min_{u \in [u_L, u_K]} a(u), & \text{si } \Lambda_{KL} \leq 0. \end{cases}$$

D'une manière simplifiée, nous écrivons $\mathcal{F}_{KL}(u_K, u_L) = \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})$ à l'instant t_{n+1} avec

$$a_{KL}^{n+1} = \begin{cases} \max_{s \in I_{KL}^{n+1}} a(s) & \text{si } \Lambda_{KL} \geq 0, \\ \min_{s \in I_{KL}^{n+1}} a(s) & \text{si } \Lambda_{KL} < 0, \end{cases} \quad \text{et} \quad I_{KL}^{n+1} = \begin{cases} [u_K, u_L], & \text{si } u_K \leq u_L, \\ [u_L, u_K], & \text{si } u_L \leq u_K. \end{cases}$$

Par conséquent, nous obtenons l'approximation suivante du terme de diffusion :

$$\begin{aligned} -\frac{1}{\Delta t} \sum_{\sigma \subset \partial \omega_K} \int_{t_n}^{t_{n+1}} \int_{\sigma} \Lambda(\mathbf{x}) a(u) \nabla u \cdot \eta_{\sigma} d\sigma(\mathbf{x}) dt &\approx \sum_{T \in \mathcal{T}} \sum_{\sigma \in \partial \omega_K \cap T} \Lambda_{KL}^T a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) \\ &= \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}). \end{aligned}$$

Le terme de convection La discrétisation du terme de convection ressemble à celle du terme de diffusion ; en effet, la fonction χ s'écrit comme le produit de deux fonctions μ et a . La fonction a est approchée en utilisant un schéma de Godunov, et ensuite nous utilisons un schéma décentré amont pour approcher le flux $\mu(u) \Lambda(\mathbf{x}) \nabla v \cdot \eta_{\sigma}$ à travers l'interface σ , ce qui nécessite l'introduction d'une fonction flux numérique vérifiant les propriétés (i)–(v) de la section 1.2.1. Nous donnons deux exemples sur l'approximation μ_{KL}^{n+1} de la fonction μ de telle sorte que la fonction flux numérique vérifie (i)–(v). Le premier exemple consiste à prendre le schéma d'Engquist-Osher et le deuxième consiste à prendre le schéma de Godunov. En particulier, nous avons :

$$\bullet \quad \mu_{KL}^{n+1} = \begin{cases} \mu_{\downarrow}(u_K^{n+1}) + \mu_{\uparrow}(u_L^{n+1}), & \text{si } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) \geq 0, \\ \mu_{\uparrow}(u_K^{n+1}) + \mu_{\downarrow}(u_L^{n+1}), & \text{si } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) < 0. \end{cases}$$

Les fonctions μ_{\uparrow} et μ_{\downarrow} sont données par :

$$\mu_{\uparrow}(z) := \int_0^z (\mu'(s))^+ ds, \quad \mu_{\downarrow}(z) := - \int_0^z (\mu'(s))^- ds.$$

$$\bullet \quad \mu_{KL}^{n+1} = \begin{cases} \max_{[u_K^{n+1}, u_L^{n+1}]} \mu(u), & \text{si } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) \geq 0, \\ \min_{[u_L^{n+1}, u_K^{n+1}]} \mu(u), & \text{si } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) \geq 0, \\ \max_{[u_L^{n+1}, u_K^{n+1}]} \mu(u), & \text{si } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) < 0, \\ \min_{[u_K^{n+1}, u_L^{n+1}]} \mu(u), & \text{si } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) < 0. \end{cases}$$

Par conséquent, nous obtenons l'approximation suivante pour de terme de convection :

$$\frac{1}{\Delta t} \sum_{\sigma \subset \partial \omega_K} \int_{t_n}^{t_{n+1}} \int_{\sigma} \mu(u) \Lambda(\mathbf{x}) a(u) \nabla v \cdot \eta_{\sigma} d\sigma(\mathbf{x}) dt \approx - \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}).$$

Discrétisation de la deuxième équation du système (1.22) Dans le but de discrétiser la deuxième équation du système (1.22), nous sommes amenés à introduire les fonctions suivantes $\eta(v)$, $p(v)$,

$\Gamma(v)$ et $\phi(v)$ définies par :

$$\begin{aligned} \eta(v) &= \max(0, \min(v, 1)), \quad p(v) = \int_1^v \frac{1}{\eta(s)} ds = \begin{cases} \ln(v) & \text{if } v \in (0, 1), \\ v - 1 & \text{if } v \geq 1, \end{cases} \\ \Gamma(v) &= \int_1^v p(s) ds = \begin{cases} v \ln(v) - v + 1 & \text{if } v \in [0, 1), \\ \frac{(v-1)^2}{2} & \text{if } v \geq 1, \end{cases} \\ \phi(v) &= \int_0^v \frac{1}{\sqrt{\eta(s)}} ds = \begin{cases} \frac{\sqrt{v}-1}{2} & \text{if } v \in [0, 1), \\ v - 1 & \text{if } v \geq 1. \end{cases} \end{aligned}$$

En utilisant ces fonctions, nous obtenons une approximation pour le terme de diffusion similaire à celle obtenue pour la première équation.

En résumé, le schéma CVFE non linéaire proposé pour la discrétisation du problème (1.22)–(1.24) consiste à chercher $U = (u_K^{n+1})_{K \in \mathcal{V}, n \in [0, \dots, N+1]}$ et $V = (v_K^{n+1})_{K \in \mathcal{V}, n \in [0, \dots, N+1]}$ solution de

$$u_{\mathcal{M}}^0(\mathbf{x}) = u_K^0 = \frac{1}{m_K} \int_{\omega_K} u_0(\mathbf{y}) d\mathbf{y}, \quad v_{\mathcal{M}}^0(\mathbf{x}) = v_K^0 = \frac{1}{m_K} \int_{\omega_K} v_0(\mathbf{y}) d\mathbf{y}. \quad (1.26)$$

et pour tout $n \in \{0, \dots, N\}$

$$\begin{aligned} m_K \frac{u_K^{n+1} - u_K^n}{\Delta t} + \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) \\ - \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}) = m_K f(u_K^{n+1}), \\ m_K \frac{v_K^{n+1} - v_K^n}{\Delta t} + \sum_{\sigma_{KL} \in \mathcal{E}_K} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1})) = m_K (\alpha u_K^n - \beta v_K^{n+1}). \end{aligned} \quad (1.27)$$

Conservativité du schéma Par construction, le schéma (1.26)–(1.27) est un schéma conservatif ; en effet, en notant $F_{KL}^{n+1} = \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) - \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1})$ et $\Phi_{KL}^{n+1} = D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))$, nous avons $a_{KL}^{n+1} = a_{LK}^{n+1}$, $\eta_{KL}^{n+1} = \eta_{LK}^{n+1}$ et $\mu_{KL}^{n+1} = \mu_{LK}^{n+1}$. Par conséquent, nous obtenons la forme conservative locale

$$\begin{cases} F_{KL}^{n+1} + F_{LK}^{n+1} = 0 = \Phi_{KL}^{n+1} + \Phi_{LK}^{n+1}, & \text{pour tout } \sigma_{KL} \in \mathcal{E}, \\ m_K \frac{u_K^{n+1} - u_K^n}{\Delta t} + \sum_{\sigma_{KL} \in \mathcal{E}_K} F_{KL}^{n+1} = f(u_K^{n+1}) m_K, & \text{pour tout } K \in \mathcal{V}, \\ m_K \frac{v_K^{n+1} - v_K^n}{\Delta t} + \sum_{\sigma_{KL} \in \mathcal{E}_K} \Phi_{KL}^{n+1} = g(u_K^n, v_K^{n+1}) m_K, & \text{pour tout } K \in \mathcal{V}. \end{cases}$$

Afin d'établir l'existence d'une solution et la convergence du schéma non linéaire (1.26)–(1.27), nous utilisons des propriétés discrètes ainsi que des estimations *a priori*. Toutes ces propriétés seront détaillées dans le chapitre. Nous commençons par la propriété suivante

$$\int_{\Omega} \mathbf{T}(\mathbf{x}) \nabla u_{\mathcal{T}} \cdot \nabla v_{\mathcal{T}} d\mathbf{x} = \sum_{\sigma_{KL} \in \mathcal{E}} a_{KL} (u_K - u_L) (v_K - v_L), \quad \mathbf{T}(\mathbf{x}) = \Lambda(\mathbf{x}) \text{ ou } D(\mathbf{x}).$$

En effet, notons que nous avons $\mathbf{T}_{KK} = -\sum_{L \neq K} \mathbf{T}_{KL}$. Par conséquent,

$$\begin{aligned} \int_{\Omega} \mathbf{T}(\mathbf{x}) \nabla u_{\mathcal{T}} \cdot \nabla v_{\mathcal{T}} d\mathbf{x} &= \sum_{T \in \mathcal{T}} \int_T \mathbf{T}(\mathbf{x}) \nabla u_{\mathcal{T}} \cdot \nabla v_{\mathcal{T}} d\mathbf{x} \\ &= - \sum_{K \in \mathcal{V}} \sum_{L \in \mathcal{V}, L \neq K} \mathbf{T}_{KL} u_K v_L - \sum_{K \in \mathcal{V}} \mathbf{T}_{KK} u_K v_K = \sum_{K \in \mathcal{V}} \sum_{L \in \mathcal{V}, L \neq K} \mathbf{T}_{KL} (v_K - v_L) u_K \\ &= \sum_{K \in \mathcal{V}} \sum_{L \in \mathcal{V}, L \neq K} \mathbf{T}_{KL} (v_K - v_L) u_K = \sum_{\sigma_{KL} \in \mathcal{E}} \mathbf{T}_{KL} (v_K - v_L) (u_K - u_L). \end{aligned}$$

Nous donnons les deux lemmes suivants dont la démonstration est une conséquence directe de la propriété précédente ainsi que de la définition des fonctions ξ et ϕ .

Proposition 1.9. *Notons $\xi_{\mathcal{T}, \Delta t}$ (resp. $\phi_{\mathcal{T}, \Delta t}$) l'unique fonction de $\mathcal{H}_{\mathcal{T}, \Delta t}$ avec des valeurs nodales $(\xi(u_K^{n+1})) \in \mathbb{R}^{(N+1)\#\mathcal{V}}$ (resp. $(\phi(v_K^{n+1})) \in \mathbb{R}^{(N+1)\#\mathcal{V}}$). Nous avons alors*

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 \\ \geq \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} (\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2 = \iint_{Q_{t_f}} \Lambda \nabla \xi_{\mathcal{T}, \Delta t} \cdot \nabla \xi_{\mathcal{T}, \Delta t} d\mathbf{x} dt, \end{aligned}$$

et

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2 \\ \geq \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} (\phi(v_K^{n+1}) - \phi(v_L^{n+1}))^2 = \iint_{Q_{t_f}} D \nabla \phi_{\mathcal{T}, \Delta t} \cdot \nabla \phi_{\mathcal{T}, \Delta t} d\mathbf{x} dt. \end{aligned}$$

Proposition 1.10. *Il existe une constante $C_1 > 0$ qui dépend seulement de Λ et $\theta_{\mathcal{T}}$ telle que*

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |\Lambda_{KL}| a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 &\leq C_1 \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2. \\ \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |D_{KL}| \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2 \\ &\leq C_1 \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2. \end{aligned}$$

Estimations discrètes, existence et convergence du schéma

Nous donnons les estimations *a priori* sur les solutions discrètes du schéma non linéaire (1.26)–(1.27). Nous présentons tout d'abord le principe de maximum discret :

Proposition 1.11. *Soit $(u_K^{n+1}, v_K^{n+1})_{K \in \mathcal{V}, n \in \{0, \dots, N\}}$ une solution du schéma non linéaire (1.26)–(1.27). Alors, pour tout $K \in \mathcal{V}_h$, et pour tout $n \in \{0, \dots, N+1\}$, nous avons $0 \leq u_K^n \leq 1$ et $v_K^n \geq 0$.*

La preuve de cette proposition est classique. Nous prenons un volume de contrôle dual ω_K tel que $u_K^{n+1} = \min \{u_L^{n+1}\}_{L \in \mathcal{V}}$, nous multiplions la première équation du système (1.27) par $-(u_K^{n+1})^-$ et nous faisons la sommation sur $K \in \mathcal{V}$. Nous obtenons la positivité de u_K^{n+1} en utilisant le prolongement par continuité des fonctions μ et f pour $u \leq 0$ et en remarquant que $a_{KL}^{n+1} = 0$ quand $\Lambda_{KL} \geq 0$, et donc $a_{KL}^{n+1} (\Lambda_{KL})^+ (u_K^{n+1} - u_L^{n+1}) (u_K^{n+1})^- \geq 0$ et par suite $\mu_{KL}^{n+1} \Lambda_{KL}^+ (v_K^{n+1} - v_L^{n+1})^- (u_K^{n+1})^- = 0$. D'une manière similaire, nous multiplions l'équation par $(u_K^{n+1} - 1)^+$ pour montrer que $u_K^{n+1} \leq 1$. La positivité de v_K^n découle directement de la positivité de u_K^n pour tout $K \in \mathcal{V}$ et tout $n \in \{0, \dots, N+1\}$.

Estimations a priori Nous avons les estimations d'énergie suivantes :

Proposition 1.12. *Pour tout $n \geq 0$, il existe une constante $C > 0$ indépendante de h et telle que*

$$\begin{aligned} \sum_{K \in \mathcal{V}} m_K \Gamma(v_K^{n+1}) + \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} (\phi(v_K^{n+1}) - \phi(v_L^{n+1}))^2 \\ \leq \sum_{K \in \mathcal{V}} m_K \Gamma(v_K^{n+1}) + \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2 \leq C. \end{aligned}$$

Proposition 1.13. *Il existe une constante $C > 0$ indépendante de h et telle que*

$$\iint_{Q_{t_f}} \Lambda(\mathbf{x}) \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) \cdot \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) \, d\mathbf{x} \, dt = \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} (v_K^{n+1} - v_L^{n+1})^2 \leq C.$$

En utilisant les propositions précédentes, nous obtenons des estimations analogues à celles données par la proposition 1.12. Spécifiquement, nous avons

Proposition 1.14. *Pour tout $n \geq 0$, il existe une constante $C > 0$ indépendante de h et telle que*

$$\begin{aligned} \sum_{K \in \mathcal{V}} m_K (u_K^{n+1})^2 + \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} (\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2 \\ \leq \sum_{K \in \mathcal{V}} m_K (u_K^{n+1})^2 + \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 \leq C. \end{aligned}$$

Dans tout ce qui suite, nous notons par $u_{\mathcal{M}, \Delta t}$ et $v_{\mathcal{M}, \Delta t}$ les éléments uniques de l'espace $\mathcal{X}_{\mathcal{M}, \Delta t}$ tels que

$$u_{\mathcal{M}, \Delta t}(\mathbf{x}_K, t) = u_K^{n+1}, \quad v_{\mathcal{M}, \Delta t}(\mathbf{x}_K, t) = v_K^{n+1}, \quad \forall K \in \mathcal{V}, \forall n \geq 0.$$

Nous montrons une estimation portant sur un raffinement amélioré de $v_{\mathcal{M}, \Delta t}$. Cette estimation prétend que la fonction $v_{\mathcal{M}, \Delta t}$ est ou bien une constante égale à zéro ou bien $v_{\mathcal{M}, \Delta t} > 0$.

Proposition 1.15. *Supposons que $\int_{\Omega} u_0(\mathbf{x}) \, d\mathbf{x} > 0$ ou $\int_{\Omega} v_0(\mathbf{x}) \, d\mathbf{x} > 0$. Alors il existe une constante r_h qui dépend des données et aussi de Δt et du maillage \mathcal{T} telle que :*

$$v_K^{n+1} \geq r_h, \quad \forall K \in \mathcal{V}, \forall n \in \{0, \dots, N\}.$$

En utilisant les propositions 1.12–1.15, nous nous appuyons sur l'argument de degré topologique pour obtenir le résultat d'existence d'une solution discrète.

Proposition 1.16 (Existence d'une solution discrète). *Soit $(u_K^n, v_K^n)_{K \in \mathcal{V}}$ un vecteur donné tel que $u_{\mathcal{M}, \Delta t}(\cdot, t_n)$ et $v_{\mathcal{M}, \Delta t}(\cdot, t_n)$ sont positives. Alors, le schéma non linéaire (1.26)–(1.26) admet au moins une solution $(u_K^{n+1}, v_K^{n+1})_{K \in \mathcal{V}}$. En outre, $u_{\mathcal{M}, \Delta t}(\cdot, t_{n+1})$ and $v_{\mathcal{M}, \Delta t}(\cdot, t_{n+1})$ sont positives.*

Le résultat principal de ce travail est le théorème suivant :

Theorem 1.17. *Soient $(\mathcal{T}_m)_{m \geq 1}$ une suite de triangulations conformes de Ω telle que $h_{\mathcal{T}_m} \rightarrow 0$ quand $m \rightarrow \infty$, et $(\Delta t_m)_{m \geq 1}$ une suite de pas de temps telle que $\Delta t_m \rightarrow 0$ quand $m \rightarrow \infty$. Pour tout $q \in [1, \infty)$, la solution discrète $(u_{\mathcal{M}_m, \Delta t_m}, v_{\mathcal{M}_m, \Delta t_m})_m$ converge faiblement dans $L^q(Q_{t_f})$ vers la solution faible du problème continue (1.22)–(1.24) quand $m \rightarrow \infty$.*

Afin d’obtenir le résultat de convergence, nous utilisons le critère de compacité de Kolmogorov en nous appuyant sur certaines estimations sur les translations en temps et en espace des solutions discrètes. Ensuite, nous identifions la solution limite avec la solution faible du problème continue (1.22)–(1.24) au sens de la définition 1.8.

La dernière partie de ce chapitre est dédiée à faire des simulations numériques en dimension 2 afin de vérifier l’efficacité du schéma proposé pour approcher les problèmes paraboliques non linéaires anisotropes modélisant la chimiotaxie.

1.2.4 Chapitre 5 : Analyse d’une équation parabolique dégénérée modélisant la chimiotaxie ou les fluides compressibles en milieu poreux

Dans ce chapitre, nous nous sommes intéressés à une équation parabolique dégénérée et non linéaire découlant, soit de la modélisation d’un fluide compressible, soit de la modélisation de la chimiotaxie. L’équation faisons l’objet d’étude dans ce chapitre est donnée, dans $Q_T = \Omega \times (0, T)$, par l’équation parabolique suivante :

$$\partial_t u - \operatorname{div}(a(u) \nabla u - f(u) \mathbf{V}) - g(u) \operatorname{div}(\mathbf{V}) + \gamma a(u) \nabla u \cdot \tilde{\mathbf{V}} = 0. \quad (1.28)$$

Cette équation découle, soit de la modélisation de la chimiotaxie-fluide (avec $\tilde{\mathbf{V}}$ est la vitesse du fluide transportant la densité cellulaire u), soit de la loi de conservation de la masse du fluide noté u (et dans ce cas $\mathbf{V} = \tilde{\mathbf{V}}$). Le domaine occupé par le fluide est noté Ω , qui est un ensemble ouvert et borné de \mathbb{R}^d , $d = 2, 3$. T représente un instant fixé dans le temps, pour lequel nous posons $\Sigma_T = \partial\Omega \times (0, T)$. À l’équation (1.28), nous ajoutons la condition initiale donnée par :

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad \text{dans } \Omega, \quad (1.29)$$

pour tout $\mathbf{x} \in \Omega$, et les conditions aux limites sur le bord Σ_T données par :

$$u(\mathbf{x}, t) = 0, \quad \text{sur } \Sigma_T, \quad (1.30)$$

Nous introduisons les hypothèses classiques en milieu poreux portant sur le système : (1.28)–(1.30)

(H1) Les champs de vitesse \mathbf{V} et $\tilde{\mathbf{V}}$ sont deux fonctions mesurables dans $(L^\infty(\Omega))^d$. En outre, le champ de vitesse \mathbf{V} vérifie la condition suivante :

$$\mathbf{V} \cdot \mathbf{n} \leq 0, \quad \text{sur } \Sigma_T,$$

où \mathbf{n} est le vecteur normal unitaire à $\partial\Omega$ dirigé vers l’extérieur de Ω .

(H2) f est une fonction dérivable dans $[0, 1]$ et $g \in C^1([0, 1])$ vérifiant

$$g(0) = f(0) = 0, \quad f(1) = g(1) = 1, \quad \text{et } g'(u) > 0 \quad \forall u \in [0, 1].$$

(H3) La condition initiale u_0 satisfait $u_0 \in L^2(\Omega)$ et $0 \leq u_0(\mathbf{x}) \leq 1$ pour presque tout $\mathbf{x} \in \Omega$.

Une difficulté majeure concernant le système (1.28)–(1.30) est la possibilité de dégénérescence du terme de diffusion. Nous donnons trois hypothèses de dégénérescence de la fonction de dissipation a qui correspondent à des situations physiques distinctes :

- (H4a) $a \in \mathcal{C}^1([0, 1], \mathbb{R})$, $a(u) > 0$ pour tout $0 < u < 1$, $a(0) > 0$, $a(1) = 0$.
 De plus, il existe $a_0 > 0$, $0 < r_2 \leq 2$, $u_* < 1$, m_1 et $M_1 > 0$ tels que
 $a(u) \geq a_0$ pour tout $0 \leq u \leq u_*$,
 $m_1(1-u)^{r_2} \leq a(u) \leq M_1(1-u)^{r_2}$, pour tout $u_* \leq u \leq 1$. En outre, il existe $c_1, c_2 > 0$
 tels que : $c_1(1-u)^{-1} \leq (f(u) - g(u))^{-1} \leq c_2(1-u)^{-1}$, pour tout $u_* \leq u < 1$.
- (H4b) $a \in \mathcal{C}^1([0, 1], \mathbb{R})$, $a(u) > 0$ pour tout $0 < u < 1$, $a(0) = 0$, $a(1) > 0$,
 De plus, il existe $r_1 > 0$, m_1 et $M_1 > 0$ tels que
 $m_1 r_1 u^{r_1-1} \leq a'(u) \leq M_1 r_1 u^{r_1-1}$, pour tout $0 \leq u \leq 1$. En outre, il existe une constant
 $C > 0$ telle que : $|f(u) - g(u)| \leq Cu$ pour tout $0 \leq u \leq 1$.
- (H4c) $a \in \mathcal{C}^1([0, 1], \mathbb{R})$, $a(u) > 0$ pour tout $0 < u < 1$, $a(0) = 0$, $a(1) = 0$,
 De plus, il existe $r_1 > 0$, $r_2 > 0$, $u_* < 1$, m_1 et $M_1 > 0$ tels que
 $m_1 r_1 u^{r_1-1} \leq a'(u) \leq M_1 r_1 u^{r_1-1}$, pour tout $0 \leq u \leq u_*$,
 $-r_2 M_1 (1-u)^{r_2-1} \leq a'(u) \leq -r_2 m_1 (1-u)^{r_2-1}$, pour tout $u_* \leq u \leq 1$. En outre, il
 existe $c_1, c_2, C > 0$ tels que : $|f(u) - g(u)| \leq Cu$ pour tout $0 \leq u \leq u_*$ et $c_1(1-u)^{-1} \leq$
 $(f(u) - g(u))^{-1} \leq c_2(1-u)^{-1}$, pour tout $u_* \leq u < 1$.

Dans tout ce que suit, nous supposons que $\gamma = 1$ et que $\tilde{\mathbf{V}} = \mathbf{V}$. Nous établissons, selon les suppositions (H4a), (H4b) ou (H4c) trois résultats d'existence. Le premier résultat concerne l'existence des solutions du système (1.28)–(1.30) dans un sens classique, alors que les deux autres concernent l'existence des solutions dans un sens affaibli par rapport à la formulation classique.

Le premier résultat : les solutions faibles et classiques Nous définissons la fonction $k \in \mathcal{C}^1([0, 1])$ par :

$$\begin{aligned} k(u) &= u, & \text{if } 0 \leq u \leq u_*, \\ k'(u) &= (f(u) - g(u))^{-1} g'(u) k(u), & \text{if } u_* \leq u < 1. \end{aligned} \quad (1.31)$$

Nous notons L la primitive de la fonction k définie par $L(u) = \int_0^u k(\tau) d\tau$, pour tout $0 \leq u < 1$.

Definition 1.18. Sous les hypothèses (H1)–(H3) et (H4a), et pour une donnée initiale u_0 satisfaisant $L(u_0) \in L^1(\Omega)$, nous disons que u est une solution faible classique du système (1.28)–(1.30) si u vérifie

$$\begin{aligned} 0 &\leq u(\mathbf{x}, t) \leq 1 \text{ pour presque tout } (\mathbf{x}, t) \in \Omega \times (0, T), \\ u &\in L^2(0, T; H_0^1(\Omega)) \cap C^0([0, T]; L^2(\Omega)), \\ \partial_t u &\in L^2(0, T; H^{-1}(\Omega)), \end{aligned}$$

et telle que

$$\begin{aligned} \int_0^T \langle \partial_t u, \varphi \rangle dt + \int_{Q_T} a(u) \nabla u \cdot \nabla \varphi d\mathbf{x} dt - \int_{Q_T} f(u) \mathbf{V} \cdot \nabla \varphi d\mathbf{x} dt \\ + \int_{Q_T} g(u) \mathbf{V} \cdot \nabla \varphi d\mathbf{x} dt + \int_{Q_T} g'(u) \mathbf{V} \cdot \nabla u \varphi d\mathbf{x} dt \\ + \int_{Q_T} a(u) \mathbf{V} \cdot \nabla u \varphi d\mathbf{x} dt = 0, \quad \forall \varphi \in L^2(0, T; H_0^1(\Omega)). \end{aligned} \quad (1.32)$$

La notation $\langle \cdot, \cdot \rangle$ représente le crochet de dualité entre $H^{-1}(\Omega)$ et $H_0^1(\Omega)$.

Theorem 1.19. Sous les hypothèses (H1) – (H3) et (H4a), il existe au moins une solution faible classique du système (1.28)–(1.30) au sens de la définition 1.18.

Le deuxième résultat : les solutions faibles dégénérées Afin de définir une notion d'une solution faible adaptée avec l'hypothèse (H4b), nous introduisons les fonctions β et h définies sur \mathbb{R} par :

$$\beta(u) = u^{r-1}, \quad h(u) = \int_0^u \beta(\tau) d\tau, \quad \text{où } r = \begin{cases} r_1 + 2, & \text{si } r_1 \leq 1, \\ r_1, & \text{si } r_1 > 1. \end{cases}$$

r_1 est la même constante définie dans l'hypothèse (H4b).

Pour un paramètre fixé $\theta \geq 0$, nous considérons les fonctions β_θ et h_θ définies par :

$$\beta_\theta(u) = u^{r-1+\theta}, \quad h_\theta(u) = \int_0^u \beta_\theta(\tau) d\tau.$$

Definition 1.20. Soit $\theta \geq 7r_1 + 6 - r$, sous les hypothèses (H1)–(H3) et (H4b). Nous disons que u est une solution faible dégénérée du système (1.28)–(1.30) si

$$\begin{aligned} 0 &\leq u(\mathbf{x}, t) \leq 1 \text{ pour presque tout } (\mathbf{x}, t) \in \Omega \times (0, T), \\ h_\theta(u) &\in L^2(0, T; H_0^1(\Omega)), \\ \sqrt{\beta'(u)} \nabla u &\in (L^2(Q_T))^d, \end{aligned}$$

et telle que la fonction F définie, pour tout $\chi \in L^2(0, T; H^1(\Omega))$ par :

$$\begin{aligned} F(u, \chi) &= - \int_{Q_T} h_\theta(u) \partial_t \chi d\mathbf{x} dt - \int_\Omega h_\theta(u_0) \chi(\mathbf{x}, 0) d\mathbf{x} + \int_{Q_T} a(u) \nabla u \cdot \nabla (\beta_\theta(u) \chi) d\mathbf{x} dt \\ &\quad + \int_{Q_T} a(u) \mathbf{V} \cdot \nabla u \beta_\theta(u) \chi d\mathbf{x} dt - \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla (\beta_\theta(u) \chi) d\mathbf{x} dt \\ &\quad + \int_{Q_T} g'(u) \mathbf{V} \cdot \nabla u \beta_\theta(u) \chi d\mathbf{x} dt, \end{aligned}$$

vérifie

$$F(u, \chi) \leq 0, \quad \forall \chi \in C^1([0, T]; H_0^1(\Omega)) \text{ avec } \chi(\cdot, T) = 0 \text{ et } \chi \geq 0.$$

De plus, elle vérifie

$$\begin{aligned} \forall \varepsilon > 0, \exists Q^\varepsilon \subset Q_T \text{ tel que } \text{mes}(Q^\varepsilon) < \varepsilon, \text{ et} \\ F(u, \chi) &= 0, \quad \forall \chi \in C^1([0, T]; H_0^1(\Omega)), \text{ supp } \chi \subset ([0, T] \times \Omega) \setminus Q^\varepsilon \end{aligned}$$

Theorem 1.21. Sous les hypothèses (H1) – (H3) et (H4b), il existe au moins une solution faible dégénérée du système (1.28)–(1.30) au sens de la définition 1.20.

Le troisième résultat : les solutions faibles dégénérées Afin de définir une notion de solution faible adaptée avec l'hypothèse (H4c), nous introduisons la fonction continue $j_{\theta, \lambda}$ définie par :

$$j_{\theta, \lambda}(u) = \begin{cases} \beta_\theta(u), & \text{si } 0 \leq u \leq u_* \\ \beta_\theta(u_*) (1 - u_*)^{1 - \frac{r'}{2} - \lambda} (1 - u)^{\frac{r'}{2} - 1} + \lambda, & \text{si } u \geq u_*. \end{cases}$$

où $r' \geq \max(2, r_2)$ (r_2 est la constante définie dans (H4c)). Nous définissons encore la fonction $J_{\theta, \lambda}$ qui est la primitive de la fonction $j_{\theta, \lambda}$ et notons $j = j_{0,0}$ et $J = J_{0,0}$. Enfin, nous introduisons les fonctions μ et G définies par :

$$\begin{cases} \mu(u) = \beta(u), \\ \mu'(u) = (f(u) - g(u))^{-1} g'(u) \mu(u), \end{cases} \quad \begin{aligned} &0 \leq u \leq u_* \text{ et } \\ &u_* \leq u < 1, \end{aligned} \quad \text{et} \quad G(u) = \int_0^u \mu(y) dy.$$

Definition 1.22. Soit $\theta \geq 7r_1 + 6 - r$, $\lambda \geq 7r_2 + 6 - \frac{r'}{2}$. Sous les hypothèses (H1)–(H3) et (H4c), nous disons que u est une solution faible dégénérée du système (1.28)–(1.30) si

$$0 \leq u(\mathbf{x}, t) \leq 1 \text{ pour presque tout } (\mathbf{x}, t) \in \Omega \times (0, T),$$

$$J(u) \in L^2(0, T; H_0^1(\Omega)), \quad \mu'^{\frac{1}{2}}(u) a^{\frac{1}{2}}(u) \nabla u \in (L^2(Q_T))^d,$$

et telle que la fonction F définie par

$$\begin{aligned} F(u, \chi) = & - \int_{Q_T} J_{\theta, \lambda}(u) \partial_t \chi \, d\mathbf{x} \, dt - \int_{\Omega} J_{\theta, \lambda}(u_0(\mathbf{x})) \chi(\mathbf{x}, 0) \, d\mathbf{x} \\ & + \int_{Q_T} a(u) \nabla u \cdot \nabla (j_{\theta, \lambda}(u) \chi) \, d\mathbf{x} \, dt + \int_{Q_T} a(u) \mathbf{V} \cdot \nabla u j_{\theta, \lambda}(u) \chi \, d\mathbf{x} \, dt \\ & - \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla j_{\theta, \lambda}(u) \chi \, d\mathbf{x} \, dt + \int_{Q_T} g'(u) \mathbf{V} \cdot \nabla u j_{\theta, \lambda}(u) \chi \, d\mathbf{x} \, dt \\ & - \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla \chi j_{\theta, \lambda}(u) \, d\mathbf{x} \, dt, \end{aligned}$$

vérifie

$$F(u, \chi) \leq 0, \quad \forall \chi \in \mathcal{C}^1([0, T]; H_0^1(\Omega)) \text{ avec } \chi(\cdot, T) = 0 \text{ et } \chi \geq 0 \quad (1.33)$$

En plus,

$$\begin{aligned} & \forall \varepsilon > 0, \exists Q^\varepsilon \subset Q_T \text{ telle que } \text{mes}(Q^\varepsilon) < \varepsilon \text{ et} \\ & F(u, \chi) = 0, \quad \forall \chi \in \mathcal{C}^1([0, T]; H_0^1(\Omega)), \text{ supp } \chi \subset ([0, T] \times \Omega) \setminus Q^\varepsilon \end{aligned} \quad (1.34)$$

Theorem 1.23. Sous les hypothèses (H1) – (H3) et (H4c), il existe au moins une solution faible dégénérée du système (1.28)–(1.30) au sens de la définition 1.22.

Liste de publications

Plusieurs publications sont issues de ce travail : un article accepté pour publication dans une revue internationale :

- M. Ibrahim, M. Saad, *On the efficacy of a control volume finite element method for the capture of patterns for a volume-filling chemotaxis model*, accepté pour publication dans **Computers & Mathematics with Applications** (disponible en ligne),

un article soumis :

- C. Cancès, M. Ibrahim, M. Saad, *On a nonlinear control volume finite element method for degenerate anisotropic reaction–convection–diffusion systems modeling the chemotaxis process*,

et un proceeding pour une conférence internationale avec comité de lecture :

- C. Cancès, M. Ibrahim, M. Saad, *A nonlinear CVFE scheme for an anisotropic degenerate nonlinear Keller–Segel model*, **Journal of Mathematics in Industry**, Springer.

Pattern formation and an upstream finite-volume scheme for a volume-filling chemotaxis model

Sommaire

2.1	Introduction	29
2.2	Volume-filling chemotaxis model	31
2.3	Pattern formation	33
2.3.1	Linear stability	33
2.3.2	Formal asymptotic expansion	34
2.3.3	Turing conditions	34
2.3.4	Bifurcation	36
2.4	Analysis for a nonlinear density	37
2.4.1	Bifurcation with chemotactic sensitivity χ	38
2.4.2	Bifurcation with growth rate α	40
2.4.3	Bifurcation with death rate β	41
2.5	Finite volume approximation	41
2.5.1	Space-time discretization and discrete functions	42
2.5.2	Finite volume scheme for system (2.29)	45
2.5.3	Numerical results	47

2.1 Introduction

Reaction–diffusion systems self-organize variety of spatio-temporal patterns such as propagating circular, spiral, and periodic waves. Those spatio-temporal patterns are observed for the first time in the Belousov-Zhabotinsky reaction system [81].

A reaction–diffusion system is generally described with a set of time-evolving partial differential equations. Each equation consists of a diffusion equation coupled with a reaction term ; the

reaction term usually describes a nonlinear phenomenon observed in nature. Most typical form of reaction–diffusion system is described with a pair of reaction–diffusion equations having an activator variable u and an inhibitor variable v , as follows

$$\begin{cases} \partial_t u = D_u \Delta u + f(u, v), \\ \partial_t v = D_v \Delta v + g(u, v), \end{cases} \quad (2.1)$$

where the functions $f(u, v)$ and $g(u, v)$ denote the reaction kinetics associated with the chemicals u and v . The variables u and v are defined in a two-dimensional bounded domain (\mathbf{x}, \mathbf{y}) and in time t with, typically, zero-flux boundary conditions; D_u and D_v are the diffusion coefficients of u and v respectively. The system (2.1) with activator and inhibitor variables is a typical model in describing the pattern formation process. Indeed, the mathematician Alan Turing established in his seminal paper [76] in 1952 that, under certain conditions, chemicals can react and diffuse in such a way as to produce inhomogeneous spatial patterns of chemical concentrations. Although a diffusion process generally brings uniform distribution of a substance, the scenario of Turing presents a nonuniform pattern as a stationary state in a reaction–diffusion system.

The idea of Turing is a simple but profound one. He said that, if in the absence of diffusion (effectively $D_u = D_v = 0$), u and v tend to a linearly stable uniform steady state then, under certain conditions called **Turing conditions**, spatially inhomogeneous patterns evolve by *diffusion driven instability* also called *Turing instability* if $D_u \neq D_v$, in more particular when the inhibitor variable v rapidly diffuses more than the activator variable u does; see, e.g. [60, 61, 58]. This concept was a novel concept because the diffusion was usually considered of having a stabilizing effect.

We can find a wide variety of reaction–diffusion systems in the fields of physics, chemistry and biology [60]. One of the most popular reaction–diffusion system is the chemotaxis model.

Chemotaxis is the feature movement of a cell along a chemical concentration gradient either towards the chemical stimulus, in this case the chemical is called chemoattractant, or away from the chemical stimulus and then the chemical is called chemorepellent. The mathematical analysis of chemotaxis models shows a plenitude of spatial patterns such as the chemotaxis models applied to skin pigmentation patterns [63, 67, 80]— that lead to aggregations of one type of pigment cell into a striped spatial pattern. And other models applied to the aggregation patterns in an epidemic disease [9], tumor growth [19], angiogenesis in tumor progression [14], and many other examples.

Theoretical and mathematical modeling of chemotaxis dates to the pioneering works of Patlak in the 1950s [69] and Keller and Segel in the 1970s [51, 52]. The review article by Horstmann [47] provides a detailed introduction into the mathematics of the Keller–Segel (KS) model for chemotaxis. In its original form, this model consists of four coupled reaction-advection-diffusion equations. These can be reduced under quasi-steady-state assumptions to a model for two unknown functions u and v . The general form of the Keller–Segel model in which we are interested is given by the following system

$$\begin{cases} \partial_t u = \operatorname{div} (D(u) \nabla u - u \zeta(u, v) \nabla v) + f(u, v), & (\mathbf{x}, t) \in Q_{t_f} = \Omega \times (0, t_f), \\ \partial_t v = D_v \Delta v + g(u, v), & (\mathbf{x}, t) \in Q_{t_f} = \Omega \times (0, t_f), \end{cases} \quad (2.2)$$

with no-flux boundary conditions on $\Sigma_{t_f} := \partial\Omega \times (0, t_f)$,

$$(D(u) \nabla u - u \zeta(u, v) \nabla v) \cdot \mathbf{n} = 0, \quad \nabla v \cdot \mathbf{n} = 0, \quad (2.3)$$

where, \mathbf{n} is the unit outward normal vector at the boundary $\partial\Omega$ of the domain Ω .

In the above model, u denotes the cell density on a given open bounded connected domain $\Omega \subset \mathbb{R}^2$ for a fixed time $t_f > 0$ and v describes the concentration of the chemoattractant. The cell

dynamics derive from population kinetics and movement, the latter comprising a diffusive flux modeling undirected (random) cell migration and an advective flux with velocity dependent on the gradient of the signal, modeling the contribution of chemotaxis. $D(u)$ describes the diffusivity of the cells (sometimes also called motility) while $\zeta(u, v)$ stands to the chemotactic sensitivity. The function $f(u, v)$ describes cell growth and cell death while $g(u, v)$ describes production and degradation of the chemoattractant. A key property of the above equations is their ability to give rise to spatial pattern formation when the chemical signal v acts as an *auto-attractant*, that is, when cells both produce and migrate up gradients of the chemical signal.

The Keller–Segel model (2.2)–(2.3) is studied by many authors, for example we can cite the review articles of Horstmann [46, 47] and the paper of Bendahmane *et al.* [8].

For a particular choice of the diffusivity of the cells $D(u)$, the chemotaxis model (2.2) transforms into a system describing the volume-filling effect. In the volume filling effect, the particles are assumed of having a finite volume and the cells cannot move into regions that are already filled by other cells. Furthermore, in the volume filling approach we consider that the cell jumped into an available space with a given squeezing probability of a cell finding space at its neighboring location see, for example [66].

In [66], Painter and Hillen introduced a squeezing probability of a cell finding space which decreases linearly with the cell density at that site; indeed, it corresponds to the interpretation that cells are solid blocks and the probability of finding space is proportional to the number of occupants (see, for example [72, 71]). But since the cells are not solid blocks and they are elastic, Wang and Hillen introduced in their paper [78] a nonlinear squeezing probability, which is greater than a linear distribution.

From a numerical point of view, we mention that Bendahmane *et al.* [7] analyzed a finite volume method for the Keller–Segel model (2.2)–(2.3). In this chapter, we adopt the same squeezing probability as in [78] and study the pattern formation for the corresponding volume-filling chemotaxis model using the linear stability analysis and finally we adopt the similar techniques used in [7] to numerically investigate the spatio-temporal patterns in a two-dimensional bounded domain.

The organization of this chapter is as follows. In section 2.2, we define and give the assumptions for the squeezing probability and for the functions f and g that make in evidence global and bounded solutions for the volume-filling chemotaxis model. In section 2.3, we establish Turing conditions for the generation of spatial patterns for the volume-filling chemotaxis model with zeros-flux boundary conditions by performing the standard linear stability analysis and by using a formal asymptotic expansion to linearize the nonlinear diffusion terms. In section 2.4, we take particular functions into the chemotaxis model, apply the pattern formation analysis, and then determine the range of wave numbers for each bifurcation parameter into the chemotaxis model. In section 2.5, we introduce a finite volume scheme to discretize the volume-filling chemotaxis model and to numerically investigate the spatio-temporal patterns in a two-dimensional bounded domain. We end up with some numerical experiments to capture the generation of spatial patterns for the volume-filling chemotaxis model.

2.2 Volume-filling chemotaxis model

In this section, we consider the volume-filling chemotaxis model (2.2)–(2.3) in which we introduce the squeezing probability q of a cell finding space at its neighboring location. q is a non-increasing function with respect to the cell density u and in addition, it is a nonlinear function and equals to zero when the maximum number of cells exceeds a certain number \bar{u} . The number \bar{u} represents the total number of cells that can be accommodated at any site.

We take the squeezing probability q introduced in [78]

$$q(u) = \begin{cases} 1 - \left(\frac{u}{\bar{u}}\right)^\gamma, & 0 \leq u \leq \bar{u}, \\ 0, & u > \bar{u}, \end{cases} \quad (2.4)$$

with, $\gamma \geq 1$ is a nonnegative exponent representing the rigidity of the cells.

In what follows, the diffusion term and the convection term are given by

$$D(u) = d_1 (q(u) - q'(u)u), \quad \zeta(u, v) = q(u)\chi(v),$$

where q is the squeezing probability of a cell finding space into its neighboring location, it is defined in (2.4) and χ is a v -depending continuous function.

Substituting these coefficients into system (2.2), we obtain the final form of the volume-filling chemotaxis model.

We denote by Ω an open bounded domain in \mathbb{R}^d , $d = 2, 3$. The final form of the volume-filling chemotaxis model is given by

$$\begin{cases} \partial_t u = \operatorname{div} (d_1 (q(u) - q'(u)u) \nabla u - q(u)u\chi(v) \nabla v) + f(u, v), & (\mathbf{x}, t) \in Q_{t_f}, \\ \partial_t v = d_2 \Delta v + g(u, v), & (\mathbf{x}, t) \in Q_{t_f}. \end{cases} \quad (2.5)$$

With the following zero flux boundary conditions

$$\begin{aligned} (d_1 (q(u) - q'(u)u) \nabla u) \cdot \mathbf{n} - q(u)u\chi(v) \nabla v \cdot \mathbf{n} &= 0, \\ \nabla v \cdot \mathbf{n} &= 0, \end{aligned} \quad (2.6)$$

where \mathbf{n} is the unit normal vector at the boundary $\partial\Omega$ and outward to Ω . We assume that the initial conditions are positive

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}) \geq 0, \quad v(\mathbf{x}, 0) = v_0(\mathbf{x}) \geq 0, \quad \text{for all } \mathbf{x} \in \Omega.$$

Furthermore, we introduce the assumptions on f , g and q for the global existence of solutions for the system (2.5)

(A1) d_1 and d_2 are two nonnegative constants.

(A2) The chemosensitivity $\chi \in \mathcal{C}^2(\mathbb{R}, \mathbb{R})$ such that $\chi(v) > 0$.

(A3) The squeezing probability q is a nonincreasing and concave function lying into $\mathcal{C}^3(\mathbb{R}, \mathbb{R})$ and satisfying the following condition : there exists a critical number \bar{u} such that $q(0) = 1$, $q(\bar{u}) = 0$, $0 < q(u) < 1$ for $u \in (0, \bar{u})$ and $q(u) = 0$ for all $u > \bar{u}$.

(A4) f is a function lying into $\mathcal{C}^2(\mathbb{R} \times \mathbb{R})$ such that $f(0, v) \geq 0$ for all $v \geq 0$. In addition, there exists a constant $0 < u_c < \bar{u}$ such that $f(u_c, v) = 0$ and $f(u, v) < 0$ for all $u > u_c$ and $v \geq 0$.

(A5) g is a function lying into $\mathcal{C}^2(\mathbb{R} \times \mathbb{R})$ such that $g(u, 0) \geq 0$ for all $u \geq 0$. Furthermore, there exists a constant $\bar{v} > 0$ such that $g(u, \bar{v}) < 0$ for all $0 \leq u \leq \bar{u}$.

The existence of global bounded classical solutions of system (2.5) object of the boundary conditions (2.6) has been established in [78] using the maximum principle and the Amman theory [4].

2.3 Pattern formation

By definition, the patterns are the solutions of a reaction–diffusion system which are stable in time and stationary inhomogeneous in space, while the pattern formation in mathematics refers to the process that, by changing a bifurcation parameter, the spatially homogeneous steady states lose stability to spatially inhomogeneous perturbations, and stable inhomogeneous solutions arise.

In this section, we are interested in the pattern formation for system (2.2). We derive here necessary conditions for diffusion-driven instability (also called **Turing instability**) of the steady state and the initiation of spatial pattern for the above-mentioned system. The approach followed here is very routine. We look for the spatial homogeneous steady states by setting the kinetics on the right-hand side of equation (2.5) to be equal to zero,

$$f(u, v) = 0, \quad g(u, v) = 0. \quad (2.7)$$

We suppose that (u_s, v_s) is the relevant homogeneous steady state of system (2.5), that is mean that (u_s, v_s) is a nonnegative solution of equation (2.7). Since we are concerned with Turing instability, we are interested in linear instability of this steady state that is solely spatially dependent. So, in the absence of any spatial variation, the homogeneous steady state must be linearly stable : we first determine the conditions for this to hold.

With no spatial variation u and v satisfy

$$\partial_t u = f(u, v), \quad \partial_t v = g(u, v). \quad (2.8)$$

2.3.1 Linear stability

Linearizing about the steady state (u_s, v_s) is exactly the same technique used by Murray [60, 61], we set

$$\mathbf{w} = \begin{pmatrix} u - u_s \\ v - v_s \end{pmatrix}, \quad (2.9)$$

and equation (2.8) becomes for $|\mathbf{w}|$ small enough,

$$\partial_t \mathbf{w} = A \mathbf{w}, \quad A = \begin{pmatrix} f_u & f_v \\ g_u & g_v \end{pmatrix}_{u_s, v_s}, \quad (2.10)$$

where A is the stability matrix (i.e. the Jacobian matrix of system (2.8) computed at the steady state (u_s, v_s)). From now on, we take the derivatives of f and g to be evaluated at the steady state (u_s, v_s) unless stated otherwise. We look for solutions in the form

$$\mathbf{w}(t) = \sum_{\lambda} c_{\lambda} e^{\lambda t}, \quad (2.11)$$

where λ is an eigenvalue of the matrix A . The steady state $\mathbf{w} = 0$ is linearly stable if $\text{Re } \lambda < 0$ since in this case the perturbation \mathbf{w} tends to zero as $t \rightarrow \infty$. Substituting equation (2.11) into equation (2.10), one can determine the eigenvalues λ as the solutions of the characteristic polynomial

$$|A - \lambda \mathbf{I}_2| = \begin{vmatrix} f_u - \lambda & f_v \\ g_u & g_v - \lambda \end{vmatrix} = 0,$$

where \mathbf{I}_2 is the identity matrix in \mathbb{R}^2 . Hence, the eigenvalues are given by

$$\lambda_1, \lambda_2 = \frac{f_u + g_v \pm \sqrt{(f_u + g_v)^2 - 4f_v g_u}}{2}. \quad (2.12)$$

Therefore, the conditions for which the steady state (u_s, v_s) is linearly stable, i.e. $\text{Re } \lambda < 0$ where λ is the eigenvalue of problem (2.10), are guaranteed if

$$\text{tr}(A) = f_u + g_v < 0, \quad \det(A) = f_u g_v - f_v g_u > 0. \quad (2.13)$$

Whether the kinetics f and g depend on some parameters, then (u_s, v_s) are functions of these parameters, thus the above inequalities impose certain conditions on these parameters.

2.3.2 Formal asymptotic expansion

Now, we consider the full chemotaxis model (2.5). We examine a small perturbation about the homogeneous steady state (u_s, v_s) in the form

$$u = u_s + \varepsilon \tilde{u}(\mathbf{x}, t), \quad v = v_s + \varepsilon \tilde{v}(\mathbf{x}, t), \quad (2.14)$$

where $\varepsilon \ll 1$ is a small parameter strictly positive.

Since f is a function of class $C^2(\mathbb{R} \times \mathbb{R})$, then f can be written using an asymptotic expansion around (u_s, v_s) as follows

$$f(u_s + \varepsilon \tilde{u}, v_s + \varepsilon \tilde{v}) = f(u_s, v_s) + \varepsilon \tilde{u} f_u(u_s, v_s) + \varepsilon \tilde{v} f_v(u_s, v_s) + \dots$$

Performing the same asymptotic expansion for the following functions g , q and χ around (u_s, v_s) , u_s , and v_s respectively, and taking into account equation (2.7), then system (2.5) transforms into the following system

$$\begin{cases} \varepsilon \partial_t \tilde{u} = \varepsilon \text{div} \left(d_1 (q(u_s + \varepsilon \tilde{u}) - q'(u_s + \varepsilon \tilde{u})(u_s + \varepsilon \tilde{u})) \nabla \tilde{u} \right. \\ \quad \left. - \varepsilon (u_s + \varepsilon \tilde{u}) \chi(v_s + \varepsilon \tilde{v}) q(u_s + \varepsilon \tilde{u}) \nabla \tilde{v} \right) \\ \quad + \varepsilon \tilde{u} f_u(u_s, v_s) + \varepsilon \tilde{v} f_v(u_s, v_s) + \dots \\ \varepsilon \partial_t \tilde{v} = \varepsilon d_2 \Delta \tilde{v} + \varepsilon \tilde{u} g_u(u_s, v_s) + \varepsilon \tilde{v} g_v(u_s, v_s) + \dots \end{cases} \quad (2.15)$$

Equating first-order terms with respect to ε , neglecting higher-order terms, and dropping the tilde for the sake of brevity, we obtain the following linearized system

$$\begin{cases} \partial_t u = d_1 (q(u_s) - q'(u_s) u_s) \Delta u - u_s q(u_s) \chi(v_s) \Delta v + u f_u + v f_v, \\ \partial_t v = d_2 \Delta v + u g_u + v g_v. \end{cases} \quad (2.16)$$

An importance of the asymptotic expansion is to transform the nonlinear coefficients of the diffusion and convection terms into constant coefficients which are nothing other than the old coefficients evaluated at the steady state (u_s, v_s) .

For simplicity, we denote by $\xi = q(u_s) - q'(u_s) u_s$ and $\psi = -q(u_s) \chi(v_s) u_s$ these coefficients. The quantity v_s in $\chi(v_s)$ will be often abbreviated for notational convenience unless stated otherwise, i.e. $\chi = \chi(v_s)$. Taking into accounts the assumptions for which there exists a nonnegative solution, then one can deduce that $\xi > 0$ and $\psi < 0$.

2.3.3 Turing conditions

In what follows, we assume that the domain is a two dimensional bounded domain. We perform the standard argument (see e.g. [61, chapter II]). We linearize again about the steady state, which with equation (2.9) is $\mathbf{w} = 0$, to get

$$\partial_t \mathbf{w} = A \mathbf{w} + \mathcal{D} \Delta \mathbf{w}, \quad \mathcal{D} = \begin{pmatrix} d_1 \xi & \psi \\ 0 & d_2 \end{pmatrix}. \quad (2.17)$$

To solve this system of equations under the Neumann boundary conditions (due to conditions (2.6) mentioned in [section 2.2](#) and to the asymptotic expansion), we define the time-independent solution $\mathbf{W}(\mathbf{x})$ of the spatial eigenvalue problem defined by

$$-\Delta \mathbf{W} = k^2 \mathbf{W} \quad \text{in } \Omega, \quad \nabla \mathbf{W} \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega, \quad (2.18)$$

where k is the eigenvalue. For example, assume that the domain is one-dimensional, say, $0 \leq \mathbf{x} \leq L$. Then, system (2.18) is summarized by the following equation

$$\mathbf{W}''(\mathbf{x}) + k^2 \mathbf{W}(\mathbf{x}) = 0, \quad \mathbf{W}'_r(0) = \mathbf{W}'_l(L) = 0, \quad \forall 0 \leq \mathbf{x} \leq L.$$

It is easy to see that $\mathbf{W}(\mathbf{x}) = \sum_{n \in \mathbb{Z}} c_n \cos(n\pi \mathbf{x}/L)$ is a nontrivial solution where c_n is a constant different from zero. The eigenvalue in this case is $k = n\pi/L$, it is also called the wavenumber. From now on, we shall refer to k in this context as the wavenumber. With finite domains, there exists a discrete set of possible wave numbers since n is an integer.

Since the problem is linear and by the superposition principle, the linear analysis consists of looking for solutions $\mathbf{W}(\mathbf{x}, t)$ to system (2.17) in the following form

$$\mathbf{W}(\mathbf{x}, t) = \sum_k c_k e^{\lambda t} \mathbf{W}_k(\mathbf{x}), \quad (2.19)$$

where for every wavenumber k , \mathbf{W}_k is the eigenfunction of problem (2.18), c_k are constants determined by a Fourier expansion of the initial conditions in terms of $\mathbf{W}_k(\mathbf{x})$, λ is the eigenvalue which determines temporal growth and it depends on the wavenumber k .

Substituting form (2.19) into equation (2.17) with (2.18), and canceling $e^{\lambda t}$, we get, for each wavenumber k

$$\lambda \mathbf{W}_k = A \mathbf{W}_k - \mathcal{D} k^2 \mathbf{W}_k.$$

We take the nontrivial solutions for \mathbf{W}_k so the λ are nothing other than the eigenvalues of the matrix $A - \mathcal{D} k^2$, which are determined by the roots of the characteristic polynomial

$$|\lambda \mathbf{I}_2 - A + \mathcal{D} k^2| = 0, \quad \mathbf{I}_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

Consequently, the dispersion relation associated with system (2.16) is determined by the following equation

$$\lambda^2 + a(k^2) \lambda + b(k^2) = 0, \quad (2.20)$$

where $a = a(k^2)$ and $b = b(k^2)$ are two functions given by

$$\begin{aligned} a(k^2) &= (\xi d_1 + d_2) k^2 - (f_u + g_v), \\ b(k^2) &= \xi d_1 d_2 k^4 + (\psi g_u - d_2 f_u - \xi d_1 g_v) k^2 + f_u g_v - f_v g_u. \end{aligned}$$

The homogeneous steady state (u_s, v_s) is linearly stable if both solutions of equation (2.20) have $\text{Re } \lambda < 0$ in the absence of any spatial effects that is equivalent to $k = 0$. We have already imposed the constraints that the steady state is stable in the absence of any spatial effects; that is, $\text{Re } \lambda < 0$. The quadratic (2.20) in this case with the requirement that $\text{Re } \lambda < 0$ gives conditions (2.13). Now, for the steady state to be unstable to spatial perturbations, we require that $\text{Re } \lambda > 0$ for some $k \neq 0$. We remark that the sum of the eigenvalues ($= -a(k^2)$) is strictly negative since the function $a(k^2)$ is strictly positive due to condition (2.13). Thus, the instability can only happen if $b(k^2)$ becomes negative for some $k \neq 0$ i.e. the dispersion relation (2.20) has two roots with opposite signs. Carrying the function $b(k^2)$, then the last condition requires that

$$\xi d_1 d_2 k^4 + (\psi g_u - d_2 f_u - \xi d_1 g_v) k^2 + f_u g_v - f_v g_u < 0. \quad (2.21)$$

Since we required that the determinant $\det(A) = f_u g_v - f_v g_u > 0$ from condition (2.13), the only possibility for inequality (2.21) to hold, requires that the sum of the roots is positive

$$\psi g_u - d_2 f_u - \xi d_1 g_v < 0. \quad (2.22)$$

Referring to condition (2.21), and since $\xi d_1 d_2$ is positive, it is necessary that the discriminant of $b(k^2) = 0$ to be positive. That is

$$(\psi g_u - d_2 f_u - \xi d_1 g_v)^2 - 4\xi d_1 d_2 (f_u g_v - f_v g_u) > 0. \quad (2.23)$$

Applying equation (2.22), we obtain from this inequality that

$$\psi g_u - d_2 f_u - \xi d_1 g_v < -2\sqrt{\xi d_1 d_2 (f_u g_v - f_v g_u)}.$$

To recap, we have now obtain necessary conditions for the generation of spatial patterns by two-species reaction–diffusion mechanisms of the form (2.5)–(2.6) which correspond to the volume-filling chemotaxis model. For simplicity, we reproduce the conditions here :

$$\begin{aligned} f_u + g_v < 0, \quad f_u g_v - f_v g_u > 0, \quad \psi g_u - d_2 f_u - \xi d_1 g_v < 0, \\ (\psi g_u - d_2 f_u - \xi d_1 g_v)^2 - 4\xi d_1 d_2 (f_u g_v - f_v g_u) > 0. \end{aligned} \quad (2.24)$$

Remembering that all these derivatives are evaluated at the steady state (u_s, v_s) .

2.3.4 Bifurcation

Here, we are interested in the mathematical study of changes in the phase portrait for a dynamical system. Such changes called bifurcations help us to determine the parameters for which we can expect pattern formation. System (2.5) depends on many parameters as χ and others involved in the kinetic functions. As well these parameters play a significant role in the production of bifurcations. Once the bifurcation occurs (by changing the value of the bifurcation parameter), we get spatial patterns under the Turing conditions.

The strength of the chemotactic sensitivity χ plays a crucial role in pattern formation. Generally, there exists a critical value χ_c such that there is no pattern formation if χ is below this critical value χ_c , while pattern formation can be expected if χ is somewhere else.

To determine explicitly this critical value, we fix all the parameters in system (2.5) except for the chemotactic sensitivity χ . Effectively, we know from the previous analysis that the bifurcation occurs when $b_{\min} = b(k_{\min}^2) = 0$, which is equivalent to

$$\psi g_u - d_2 f_u - \xi d_1 g_v = -2\sqrt{\xi d_1 d_2 (f_u g_v - f_v g_u)}. \quad (2.25)$$

Substituting $\psi = -u_s q(u_s) \chi(v_s)$ into equation (2.25), then the critical chemosensitivity χ_c is given in the following expression

$$\chi_c = \frac{2\sqrt{\xi d_1 d_2 (f_u g_v - f_v g_u)} - d_2 f_u - \xi d_1 g_v}{g_u u_s q(u_s)}. \quad (2.26)$$

For this bifurcation value, the associated critical wavenumber is given by $b'(k_c^2) = 0$, that is

$$k_c^2 = \frac{d_2 f_u + \xi d_1 g_v + u_s q(u_s) \chi_c g_u}{2\xi d_1 d_2}. \quad (2.27)$$

Whenever $b(k^2)$ is negative, the equation (2.20) has a solution λ which is positive for the same range of wave numbers that make $b(k^2) < 0$. When $\chi > \chi_c$, the range of unstable wave numbers $k_1^2 < k^2 < k_2^2$ is obtained from the zeros k_1^2 et k_2^2 of equation $b(k^2) = 0$, and it is defined as

$$\begin{aligned} k_1^2 &= \frac{C - \sqrt{C^2 - 4\xi d_1 d_2 (f_u g_v - f_v g_u)}}{2\xi d_1 d_2} < k^2 \\ &< k_2^2 = \frac{C + \sqrt{C^2 - 4\xi d_1 d_2 (f_u g_v - f_v g_u)}}{2\xi d_1 d_2}, \end{aligned} \quad (2.28)$$

where $C = d_2 f_u + \xi d_1 g_v + u_s q(u_s) \chi > 0$ is nothing else the symmetric of the coefficient of k^2 in equation $b(k^2) = 0$.

From all the foregoing analysis, under conditions (2.24) and if there exists a range of wave numbers k satisfying condition (2.28), then the corresponding spatial eigenfunctions are linearly unstable and the generation of spatial patterns can be expected.

2.4 Analysis for a nonlinear density

Hereinafter, we introduce the probability function $q(u)$ which reflects the elastic property of particles

$$q(u) = 1 - \left(\frac{u}{\bar{u}}\right)^\gamma, \quad \gamma \geq 1, \quad 0 \leq u \leq \bar{u},$$

where \bar{u} is the total number of cells that can be accommodated at any site. We assume that the cell density and the chemoattractant kinetics are respectively given by

$$f(u, v) = \mu u \left(1 - \frac{u}{u_c}\right), \quad g(u, v) = \alpha u - \delta v.$$

This choice is the best choice for the cell kinetics where it represents the logistic growth with carrying capacity $0 < u_c \leq \bar{u}$ and $\mu > 0$; while the choice for chemoattractant kinetics is extensively used in literature (see, e.g. [48]) which represents the chemoattractant growth with rate α and the death with rate δ due to dilution.

Substituting the above mentioned functions into system (2.5) and assuming that the function χ is a constant function; then, we get the following system :

$$\begin{cases} \partial_t u = \operatorname{div} (D(u) \nabla u - \chi \varphi(u) \nabla v) + \mu u \left(1 - \frac{u}{u_c}\right), & \text{in } Q_{t_f} \\ \partial_t v = d_2 \Delta v + \alpha u - \beta v, & \text{in } Q_{t_f} \\ (D(u) \nabla u - \chi \varphi(u) \nabla v) \cdot \mathbf{n} = 0, \quad \nabla v \cdot \mathbf{n} = 0, & \text{in } \partial\Omega \times (0, t_f) \end{cases} \quad (2.29)$$

where, \mathbf{n} is the unit outward normal vector at the boundary $\partial\Omega$ of the domain Ω and the functions $D(u)$ et $\varphi(u)$ are given by

$$D(u) = d_1 \left(1 + (\gamma - 1) \left(\frac{u}{\bar{u}}\right)^\gamma\right), \quad \varphi(u) = u \left(1 - \left(\frac{u}{\bar{u}}\right)^\gamma\right).$$

The homogeneous steady states of system (2.29) are $(0, 0)$ and $(u_c, \alpha u_c / \beta)$. In addition, by linearization, one can determine that the steady state $(0, 0)$ is a saddle point and consequently unstable, while the homogeneous steady state $(u_c, \alpha u_c / \beta)$ is stable to the corresponding homogeneous system of (2.29). Indeed, computing the matrix stability at the steady state $(u_s, v_s) = (u_c, \alpha u_c / \beta)$, one can determine the eigenvalues corresponding to that matrix and find two negative eigenvalues $-\alpha$ et $-\beta$, thus conditions (2.13) are verified and we focus on this steady state to study pattern

formation for the system (2.29).

Let us begin by the pattern formation analysis for the system (2.29) associated to functions q , f , and g ; we apply an asymptotic expansion around the steady state $(u_c, \alpha u_c / \beta)$, one get

$$\begin{cases} \partial_t u = \operatorname{div} (D(u_c) \nabla u - \chi \varphi(u_c) \nabla v) - \mu u, & \text{in } Q_{t_f}, \\ \partial_t v = d_2 \Delta v + \alpha u - \beta v. & \text{in } Q_{t_f}. \end{cases} \quad (2.30)$$

Next, we follow the same analysis as in the previous section to end up with the following relation dispersion

$$\lambda^2 + a(k^2) \lambda + b(k^2) = 0, \quad (2.31)$$

where $a(k^2)$ and $b(k^2)$ are two functions given by

$$\begin{aligned} a(k^2) &= (D(u_c) + d_2) k^2 + (\mu + \beta), \\ b(k^2) &= D(u_c) d_2 k^4 + (\mu d_2 + \beta D(u_c) - \chi \alpha \varphi(u_c)) k^2 + \mu \beta. \end{aligned} \quad (2.32)$$

Thereby, taking into account the stability matrix corresponding to the steady state $(u_c, \alpha u_c / \beta)$, then conditions (2.24) are verified if and only if

$$\mu d_2 + \beta D(u_c) - \chi \alpha \varphi(u_c) < -2\sqrt{d_2 \mu \beta D(u_c)}. \quad (2.33)$$

The condition (2.33) is a necessary condition for the pattern formation of the system (2.30) with the specific functions q , f , and g . It is also a sufficient condition for an infinite domain while for a finite domain, there exists a range of discrete wave numbers depending in part on the boundary conditions, and for which we can expect the generation of spatial patterns. Furthermore, we have to study the influence of the bifurcation parameters to determine that range of wave numbers.

In the remainder of this section, we will investigate the influence of each of the parameters, on which depends the system (2.29), on the pattern formation for system (2.29). The main parameters we are interested in are : the squeezing exponent γ , the chemosensitivity χ , the growth rate α , and the death rate β of the chemoattractant. For each of these parameters, we determine the range of wave numbers for which there exists spatial pattern formation for the system (2.29).

2.4.1 Bifurcation with chemotactic sensitivity χ

From the previous section, we can think about the chemotactic sensitivity χ as a bifurcation parameter, we determine the bifurcation value χ_c to obtain

$$\chi_c = \frac{2\sqrt{d_2 \mu \beta D(u_c)} + \mu d_2 + \beta D(u_c)}{\alpha \varphi(u_c)}. \quad (2.34)$$

The corresponding critical wave number k_c is determined from equation (2.27) by

$$k_c^2 = \frac{-\mu d_2 - \beta D(u_c) + \chi_c \alpha \varphi(u_c)}{2d_2 D(u_c)}. \quad (2.35)$$

Whenever $b(k^2)$ is negative, equation (2.31) has a solution λ which is positive for the same range of wave numbers that make $b = b(k^2)$ negative. Furthermore, the range of unstable wave numbers $k_1^2 < k^2 < k_2^2$ is obtained from the zeros k_1^2 and k_2^2 of equation $b(k^2) = 0$, it is defined as

$$k_1^2 = \frac{S - \sqrt{S^2 - 4\mu d_2 \beta D(u_c)}}{2\xi d_1 d_2} < k^2 < \frac{S + \sqrt{S^2 - 4\mu d_2 \beta D(u_c)}}{2\xi d_1 d_2} = k_2^2,$$

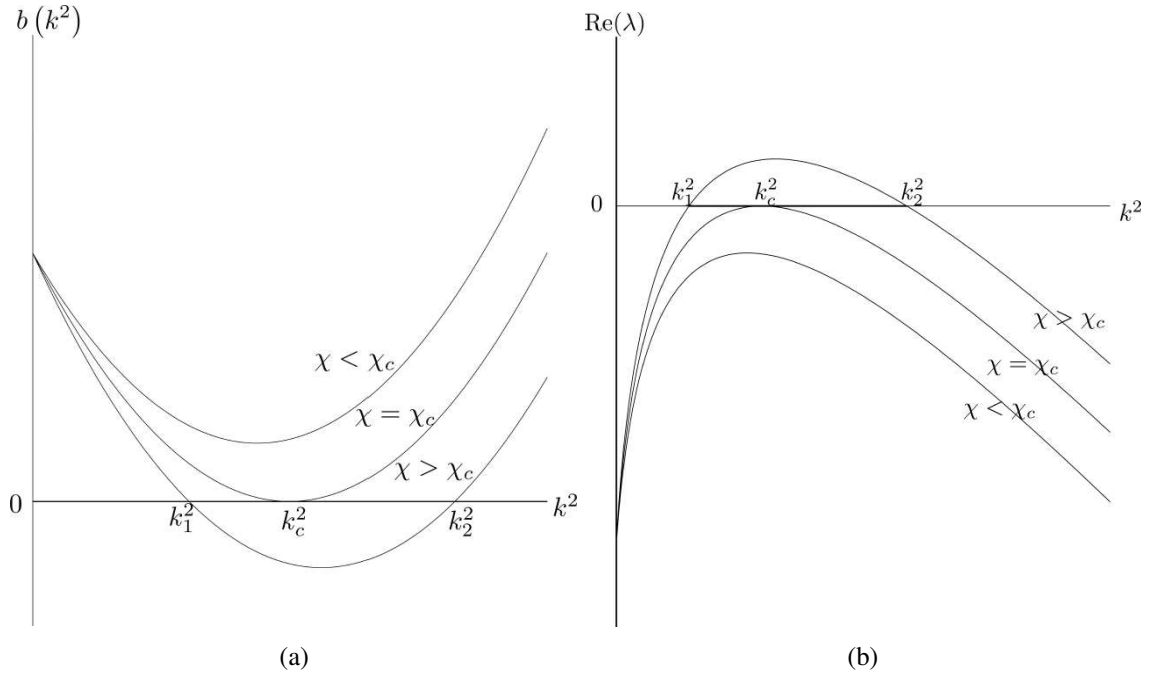


FIGURE 2.1 – (a) Plot of $b(k^2)$ as a function of k^2 defined by equation (2.32). When the chemosensitivity strength χ increases beyond the critical value χ_c , $b(k^2)$ becomes negative for a finite range of $k^2 > 0$.

(b) Plot of the real part of the eigenvalue $\lambda(k^2)$ as a function of k^2 defined by equation (2.31). When $\chi > \chi_c$, there is a range of wave numbers $k_1^2 < k^2 < k_2^2$ for which the homogeneous steady state is linearly unstable.

where $S = -\mu d_2 - \beta D(u_c) + \chi \alpha \varphi(u_c)$.

Figure 2.1a represents the variation of $b(k^2)$ as a function of k^2 and for various values of χ . Note that the critical value is (k_c, χ_c) such that $b(k^2) > 0$ if $\chi < \chi_c$, for all $k^2 > 0$, nevertheless, $b(k^2) < 0$ if $\chi > \chi_c$ and for a range of unstable wave numbers $k_1^2 < k^2 < k_2^2$.

Figure 2.1b represents the variation of the dispersion relation $\lambda(k^2)$ as a function of k^2 and for various values of χ . This graph is extremely informative in that it immediately says which eigenfunctions are linearly unstable and grow exponentially with time. Indeed, if we consider the solution \mathbf{W} given by equation (2.19), the dominant contributions as t increases are those modes for which $\text{Re } \lambda(k^2) > 0$ since all other modes tend to zero exponentially. Thus, from Fig. 2.1, we determine the range $k_1^2 < k^2 < k_2^2$, where $b(k^2) < 0$, and hence $\text{Re } \lambda(k^2) > 0$, and the eigenfunctions are unstable.

Relationship between the critical value χ_c and the exponent γ

Herein, we are interested in the relationship between the squeezing exponent γ and the critical value χ_c . A such relation is important to know the influence of the parameter γ on the dynamics of system (2.29). Indeed, the greater the critical value is smaller, the more pattern formation will be fast, here we show that the critical value χ_c as a function of the exponent γ , is nonincreasing. Hence γ plays an important role in the pattern formation for system (2.29).

We denote by M the function defined by

$$M(\gamma) = \frac{1}{1 - \left(\frac{u_c}{\bar{u}}\right)^\gamma},$$

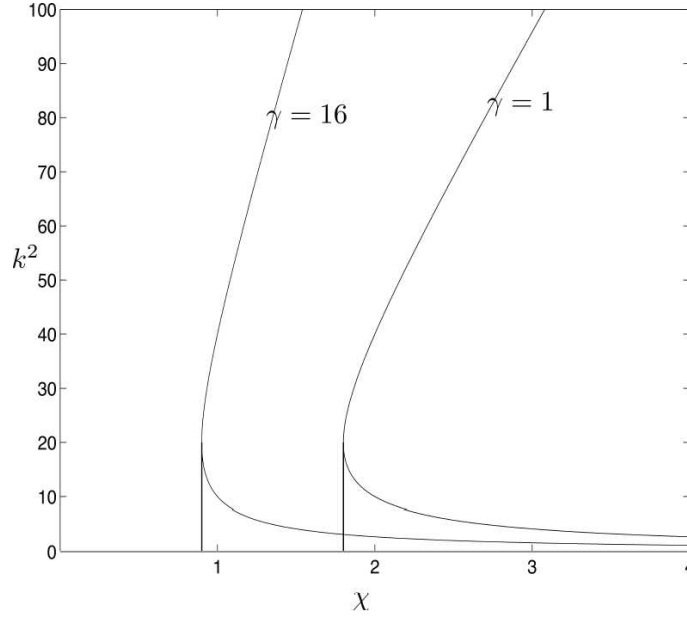


FIGURE 2.2 – Graphic representing the influence of the exponent parameter γ on the critical value χ_c . When γ increases, χ_c decreases in value. The other parameters of system (2.30) are fixed and chosen as $d_1 = 0.1$, $d_2 = 1.0$, $u_c = 2.0$, $\bar{u} = 4.0$, $\mu = 4.0$, $\alpha = 5.0$, $\beta = 10.0$.

then M is nonincreasing with respect to γ , since $u_c/\bar{u} < 1$.

Now, writing the critical value χ_c as a function of $M(\gamma)$, one can easily verify that χ_c is decreasing as a function of the exponent γ . Thereby the pattern formation is easier to form when γ is great. Figure 2.2 represents for different values of γ , the variation of the critical value (χ_c, k_c^2) with respect to γ . Note that the greater the value of γ is bigger, the more value χ_c decreases while remaining minus the value 0.9 and for the fixed parameters.

2.4.2 Bifurcation with growth rate α

In this part, we consider the growth rate α as the bifurcation parameter. We apply the same method used in section 2.4.1, i.e. we fix all the parameters in system (2.29) except the parameter α . We want to study the influence of the dynamical parameter α on the pattern formation for system (2.29).

We compute the critical value α_c for the growth rate α , we obtain

$$\alpha_c = \frac{2\sqrt{d_2\mu\beta D(u_c)} + \mu d_2 + \beta D(u_c)}{\chi\varphi(u_c)} \quad (2.36)$$

By comparing this value with the critical value χ_c in equation (2.34), one can note that there is no unstable modes if α is below the critical value α_c , whereas unstable modes are possible when α is beyond this critical value α_c .

Figure 2.3 represents the plot of the relation dispersion as a function of k^2 and for various values of α , we remark that we have consistent results to those in figure 2.1a concerning the existence of a range of wave numbers $k_1^2 < k^2 < k_2^2$ for which the eigenfunctions are unstable. Consequently, the pattern formation for system (2.29) can be expected by increasing the value of the dynamical parameter α above the critical value α_c .

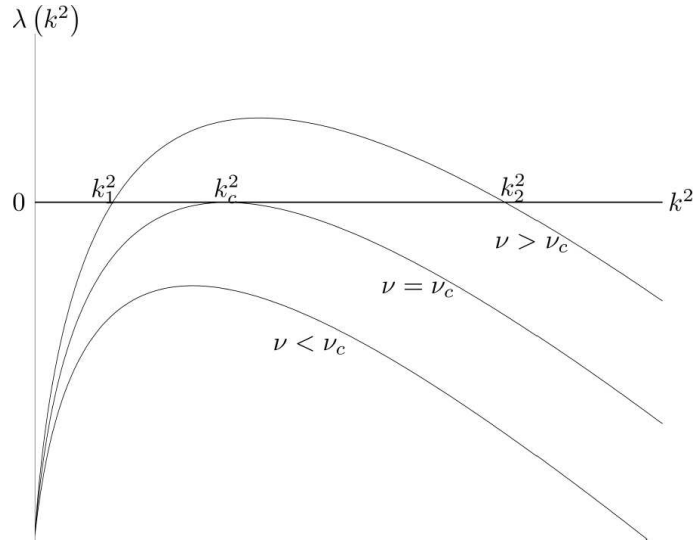


FIGURE 2.3 – Plot of the relation dispersion (2.31) as the growth rate parameter α passes through the critical value α_c . The parameters of the system (2.30) are chosen as $d_1 = 0.1$, $d_2 = 1.0$, $u_c = 2.0$, $\bar{u} = 4.0$, $\mu = 4.0$, $\alpha = 5.0$, $\beta = 10.0$, $\chi = 0.5$

2.4.3 Bifurcation with death rate β

We are interested now by the study of the influence of the death rate β on the pattern formation for system (2.29). We perform the similar linear stability as we did in the previous sections. We compute the critical value β_c for the death rate β , we obtain

$$\chi\alpha\varphi(u_c) = \mu d_2 + \beta_c D(u_c) + 2\sqrt{d_2\mu\beta_c D(u_c)} = \left(\sqrt{d_2\mu} + \sqrt{\beta_c D(u_c)}\right)^2,$$

that is

$$\beta_c = \frac{\left(\sqrt{\chi\alpha\varphi(u_c)} - \sqrt{d_2\mu}\right)^2}{D(u_c)}. \quad (2.37)$$

Taking into account the condition (2.33), one can note that there do not exist unstable modes when $\beta > \beta_c$, whereas unstable modes are possible when β is below the critical value β_c .

Fig. 2.4 represents the plot of the relation dispersion as a function of k^2 and for several values of β . We observe that the plot of the relation dispersion as β passes through the critical value β_c crosses the horizontal axis of k^2 and as result, there exists a range of unstable wave numbers $k_1^2 < k^2 < k_2^2$. Hence, the pattern formation can be expected by decreasing the value of the dynamical parameter β below the critical value β_c .

2.5 Finite volume approximation

In this section, we propose a numerical scheme to investigate the generation of spatial patterns for the volume-filling chemotaxis model (2.29) in a two dimensional space. We use the finite volume method to discretize the diffusion terms over an unstructured mesh namely admissible mesh (see [31]). Whereas, the convective term is discretized using a classical upwind finite volume scheme. The resulted scheme ensures the validity of the discrete maximum principle under some conditions on the mesh. Maintaining the discrete maximum principle is one of the important features of the finite volume scheme. The local conservativity of the numerical fluxes is an another feature of the finite volume scheme, i.e. the numerical flux is conserved from one discretization cell to its neighbor.

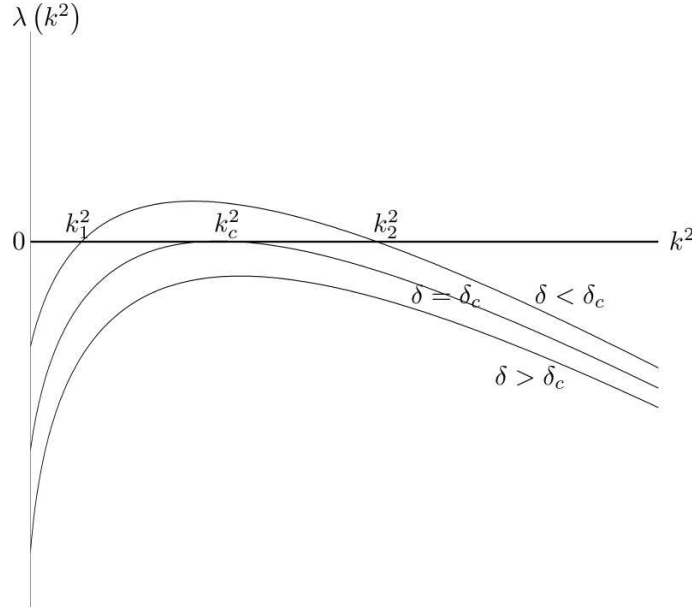


FIGURE 2.4 – Plot of the relation dispersion (2.31) as the death rate parameter β passes through the critical value β_c . For the other parameters of system (2.30), we fix $d_1 = 0.1$, $d_2 = 1.0$, $u_c = 2.0$, $\bar{u} = 4.0$, $\mu = 10.0$, $\alpha = 20$, $\chi = 0.5$, and we find that the relation dispersion passes through the bifurcation value β_c .

2.5.1 Space-time discretization and discrete functions

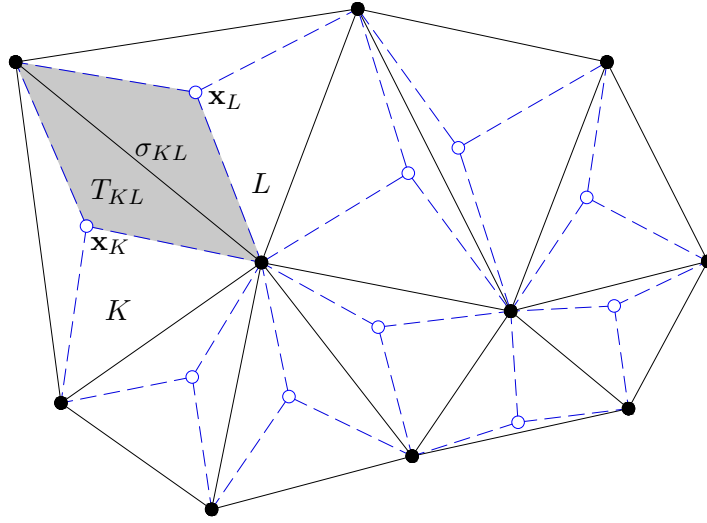
In order to discretize system (2.29), we introduce some definitions and notations. We assume that $\Omega \subset \mathbb{R}^2$, is an open bounded polygonal connected domain with boundary $\partial\Omega$.

Space discretization

Definition 2.1 (Admissible mesh). An admissible finite volume mesh of Ω is a triplet $(\mathcal{T}, \mathcal{E}, \mathcal{P})$, where \mathcal{T} is a finite family of disjoint open polygonal convex subsets of Ω called control volumes, \mathcal{E} is a finite family of subsets of Ω contained in straights of \mathbb{R}^2 with strictly positive one-dimensional measure, called the edges of the control volumes, and $\mathcal{P} = \{\mathbf{x}_K, K \in \mathcal{T}\}$ is a finite family of points of Ω , called the centers of the control volumes.

The triplet $(\mathcal{T}, \mathcal{E}, \mathcal{P})$ satisfies the following properties :

1. The closure of the union of all the control volumes is $\bar{\Omega}$, i.e. $\bar{\Omega} = \bigcup_{K \in \mathcal{T}} \bar{K}$.
2. For any $K \in \mathcal{T}$, there exists a subset \mathcal{E}_K of \mathcal{E} such that $\partial K = \bar{K} \setminus K = \bigcup_{\sigma \in \mathcal{E}_K} \bar{\sigma}$. Furthermore, $\mathcal{E} = \bigcup_{K \in \mathcal{T}} \mathcal{E}_K$.
3. For any $(K, L) \in \mathcal{T}^2$ with $K \neq L$, either the length of $\bar{K} \cap \bar{L}$ is 0 or $\bar{K} \cap \bar{L} = \bar{\sigma}$ for some $\sigma \in \mathcal{E}$, which then is denoted by σ_{KL} i.e. the interface between K and L .
4. The family $\mathcal{P} = (\mathbf{x}_K)_{K \in \mathcal{T}}$ is such that $\mathbf{x}_K \in \bar{K}$ (for all $K \in \mathcal{T}$) and, if $\sigma = \sigma_{KL}$, we assumed that $\mathbf{x}_K \neq \mathbf{x}_L$, and the straight line \mathcal{D}_{KL} going through \mathbf{x}_K and \mathbf{x}_L is orthogonal to σ_{KL} .
5. For any $\sigma \in \mathcal{E}$ such that $\sigma \subset \partial\Omega$, let K the control volume such that $\sigma \in \mathcal{E}_K$. If $\mathbf{x}_K \notin \sigma$, let $\mathcal{D}_{K,\sigma}$ be the straight line going through \mathbf{x}_K and orthogonal to σ , then the condition $\mathcal{D}_{K,\sigma} \cap \sigma \neq \emptyset$ is assumed; let $\mathbf{y}_\sigma = \mathcal{D}_{K,\sigma} \cap \sigma$.

FIGURE 2.5 – Finite volume mesh \mathcal{T} : control volumes, centers and diamonds.

In the sequel, we use the following notations. The size of the mesh \mathcal{T} is defined by :

$$h = \text{size}(\mathcal{T}) = \sup\{\text{diam}(K), K \in \mathcal{T}\},$$

where $\text{diam}(K)$ represents the least upper bound of the set of all distances between pairs of vertices of every control volume $K \in \mathcal{T}$.

For any $K \in \mathcal{T}$ and $\sigma \in \mathcal{E}$, we denote by $|K|$ the 2-dimensional Lebesgue measure of K (it is the area of K in the 2-dimensional case) and by $|\sigma|$ the length of the interface σ .

The set of neighbors of K is denoted by $\mathcal{N}(K)$, i.e. $\mathcal{N}(K) = \{L \in \mathcal{T}; \exists \sigma \in \mathcal{E}_K, \bar{\sigma} = \bar{K} \cap \bar{L}\}$; a generic neighbor of K is often denoted by L , furthermore we denote by η_{KL} and d_{KL} the unit normal vector to σ_{KL} outward to K and the distance $|\mathbf{x}_K - \mathbf{x}_L|$ respectively.

Time discretization

The time discretization is considered to be uniform (we do not impose any restriction on the time step), it is given by the sequence of discrete time $t_n = n\Delta t$ for every $n \in \mathbb{N}$ with a fixed time step Δt such that there exists an integer $N \in \mathbb{N}^*$ verifying $t_{N+1} = (N+1)\Delta t = t_f$.

Discrete space $H_{\mathcal{T}}$

Let us define the discrete finite volume space $H_{\mathcal{T}}$ of piecewise constant functions on the mesh \mathcal{T} by

$$H_{\mathcal{T}} = \{\phi : \Omega \longrightarrow \bar{\mathbb{R}}; \phi|_K \in \bar{\mathbb{R}} \text{ is constant, } \forall K \in \mathcal{T}\}.$$

Every function $u_{\mathcal{T}} \in H_{\mathcal{T}}$ is characterized by its numerical values $(u_K)_{K \in \mathcal{T}}$ such that $u_{\mathcal{T}}|_K = u_K$ for every $K \in \mathcal{T}$. In other words, given a vector $(u_K)_{K \in \mathcal{T}} \in \mathbb{R}^{\#\mathcal{T}}$; then, there exists a unique discrete function $u_{\mathcal{T}} \in H_{\mathcal{T}}$ such that

$$u_{\mathcal{T}}(\mathbf{x}) = u_{\mathcal{T}}(\mathbf{x}_K) = u_K, \quad \forall \mathbf{x} \in K, \forall K \in \mathcal{T}.$$

The discrete space $H_{\mathcal{T}}$ is a linear subspace of $L^2(\Omega)$, the usual inner scalar product becomes

$$(u_{\mathcal{T}}, v_{\mathcal{T}})_{L^2(\Omega)} = \int_{\Omega} u_{\mathcal{T}}(\mathbf{x}) v_{\mathcal{T}}(\mathbf{x}) d\mathbf{x} = \sum_{K \in \mathcal{T}} |K| u_K v_K, \quad \forall u_{\mathcal{T}}, v_{\mathcal{T}} \in H_{\mathcal{T}}.$$

Therefore, the associated norm is given by

$$\|u_{\mathcal{T}}\|_{L^2(\Omega)} = \left(\sum_{K \in \mathcal{T}} |K| |u_K|^2 \right)^{1/2}.$$

We define a discrete equivalent to the semi-norm on $H^1(\Omega)$, $|u|_{1,\Omega} = \left(\int_{\Omega} |\nabla u|^2 \right)^{1/2}$, by

$$|u_{\mathcal{T}}|_{1,\mathcal{T}} = \left(\sum_{\sigma_{KL} \in \mathcal{E}} \frac{|\sigma_{KL}|}{d_{KL}} |u_L - u_K|^2 \right)^{1/2},$$

for all $u_{\mathcal{T}} \in H_{\mathcal{T}}$.

Definition 2.2. Let $\sigma \in \mathcal{E}$ be an interface for a control volume $K \in \mathcal{T}$. A “diamond” T_{σ} constructed from the neighbor centers $\mathbf{x}_K, \mathbf{x}_L$ and the interface σ_{KL} is defined by

$$\begin{cases} T_{KL} = \bigcup_{V \in \{K,L\}} \{c\mathbf{x}_V + (1-c)\mathbf{y}, c \in [0,1[, \mathbf{y} \in \sigma_{KL}\}, & \text{if } \sigma = \sigma_{KL} \in \mathcal{E} \setminus \partial\Omega, \\ T_{K\sigma} = \{c\mathbf{x}_K + (1-c)\mathbf{y}, c \in [0,1[, \mathbf{y} \in \sigma_{KL}\}, & \text{if } \sigma \in \partial K \cap \partial\Omega. \end{cases}$$

The 2-dimensional Lebesgue measure of an interior diamond T_{KL} is $|T_{KL}| = \frac{|\sigma_{KL}|d_{KL}}{2}$.

Using definition 2.2, we define the discrete gradient over every “diamond” T_{σ} . Specifically, we have the following definition.

Definition 2.3 (Discrete gradient). The discrete gradient $\nabla_{\mathcal{T}}$ is a correspondence that maps a function $u_{\mathcal{T}} \in H_{\mathcal{T}}$ into a piecewise constant function over the “diamonds” such that

$$\nabla_{\mathcal{T}} u_{\mathcal{T}}|_{T_{\sigma}} = \begin{cases} 2 \frac{u_L - u_K}{d_{KL}} \eta_{KL}, & \text{if } \sigma = \sigma_{KL} \in \mathcal{E} \setminus \partial\Omega, \\ 0, & \text{if } \sigma \in \partial K \cap \partial\Omega. \end{cases}$$

Remark 1. From the definition 2.3 of the discrete gradient, one can deduce that

$$\|\nabla_{\mathcal{T}} u_{\mathcal{T}}\|_{L^2(\Omega)} = \left(2 \sum_{\sigma_{KL} \in \mathcal{E}} \frac{|\sigma_{KL}|}{d_{KL}} |u_L - u_K|^2 \right)^{1/2} = \sqrt{2} |u_{\mathcal{T}}|_{1,\mathcal{T}}.$$

Definition 2.4 (Discrete functions). Using the values of $(u_K^n, v_K^n), \forall K \in \mathcal{T}$ and $n \in \{0, \dots, N+1\}$, we determine the approximate finite volume solutions as piecewise constant on the control volumes and piecewise in time such that

$$\begin{cases} (u_{\mathcal{T},\Delta t}(\mathbf{x}, 0), v_{\mathcal{T},\Delta t}(\mathbf{x}, 0)) = (u_K^0, v_K^0), & \forall \mathbf{x} \in \overset{\circ}{K}, K \in \mathcal{T}, \\ (u_{\mathcal{T},\Delta t}(\mathbf{x}, t), v_{\mathcal{T},\Delta t}(\mathbf{x}, t)) = (u_K^{n+1}, v_K^{n+1}), & \forall \mathbf{x} \in \overset{\circ}{K}, K \in \mathcal{T}, \forall t \in (t_n, t_{n+1}], \end{cases}$$

where the quantity u_K^0 (resp. v_K^0) represents the mean value of the function u_0 (resp. of the function v_0). The discrete time-space of these functions is denoted by $H_{\mathcal{T},\Delta t}$.

In what follows, we use the Lipschitz nondecreasing function $A : \mathbb{R} \rightarrow \mathbb{R}$ defined by

$$A(u) = \int_0^u \Gamma(s) \, ds, \quad (2.38)$$

where $\Gamma : \mathbb{R} \rightarrow \mathbb{R}$ is the function defined by $\Gamma(u) = 1 + (\gamma - 1) \left(\frac{u}{u}\right)^{\gamma}$, so that $D(u) = d_1 \Gamma(u)$. The function A is nonlinear, we denote by $A(u_{\mathcal{T},\Delta t})$ the corresponding piecewise constant reconstruction in $H_{\mathcal{T},\Delta t}$.

2.5.2 Finite volume scheme for system (2.29)

Once the space-time discretization of the domain has been built, we can discretize the volume filling chemotaxis model (2.29). The finite volume method consists of integrating the equations of the model (contrary to the finite element method which is based on integrating the equations of the weak formulation). Applying the finite volume method, we integrate the equations of system (2.29) over the set $K \times [t_n, t_{n+1}]$ with $K \in \mathcal{T}$ and $n \in \{0, \dots, N+1\}$, further by the use of the Green-Gauss formula for the divergence operators, the finite volume consists of approximating the resulting normal fluxes across the interfaces $\{\sigma_{KL}\}_{L \in \mathcal{N}(K)}$ of the control volume K .

Time evolution terms

For the time evolution terms of system (2.29), and using the Taylor expansion, one gets

$$\begin{aligned} \int_{t_n}^{t_{n+1}} \int_K \partial_t w(\mathbf{x}, t) \, d\mathbf{x} \, dt &\approx \Delta t \int_K \partial_t w(\mathbf{x}, t_n) \, d\mathbf{x} \, dt \\ &= \int_K (w(\mathbf{x}, t_{n+1}) - w(\mathbf{x}, t_n)) \, d\mathbf{x} = |K| (w_K^{n+1} - w_K^n), \end{aligned}$$

with $w \equiv u$ or v .

Diffusion terms

We consider the diffusive term of the first equation of system (2.29), one has

$$\int_{t_n}^{t_{n+1}} \int_K \operatorname{div}(\nabla A(u)) \, d\mathbf{x} \, dt = \sum_{L \in \mathcal{N}(K)} \int_{t_n}^{t_{n+1}} \int_{\sigma_{KL}} \nabla A(u) \cdot \eta_{KL} \, d\sigma(\mathbf{x}) \, dt,$$

where, $d\sigma(\mathbf{x})$ is the Lebesgue measure on the edge σ_{KL} .

Using again Taylor's formula as well as the orthogonality condition (i.e. $\mathbf{x}_L - \mathbf{x}_K = d_{KL}\eta_{KL}$), one can deduce that the flux $\nabla A(u) \cdot \eta_{KL}$ can be approximated on the interface σ_{KL} by

$$\nabla A(u) \cdot \eta_{KL} \approx \frac{u(\mathbf{x}_L) - u(\mathbf{x}_K)}{d_{KL}}. \quad (2.39)$$

Using the flux approximation (2.39), one can deduce the following approximation of the diffusion term

$$\int_{t_n}^{t_{n+1}} \int_K \operatorname{div}(\nabla A(u)) \, d\mathbf{x} \, dt \approx \Delta t \sum_{L \in \mathcal{N}(K)} \tau_{KL} (A(u_L) - A(u_K)),$$

where, $\tau_{KL} := \frac{|\sigma_{KL}|}{d_{KL}}$ represents the transmissibility through the interface σ_{KL} . In the same manner, we treat the diffusion term of the second equation of system (2.29) to get

$$\int_{t_n}^{t_{n+1}} \int_K \operatorname{div}(\nabla v) \, d\mathbf{x} \, dt \approx \Delta t \sum_{L \in \mathcal{N}(K)} \tau_{KL} (v_L - v_K).$$

Convection terms

Let us consider the convective term of the first equation of system (2.29). For the simplicity, we denote by Ψ the function defined by $\Psi(u) := u q(u) \chi$. We note that the discretization of the convective term is slightly different from that one of the diffusion term, here we have to approximate $\Psi(u) \nabla v \cdot \eta_{KL}$ at the interface σ_{KL} and at time t_{n+1} . The classical choice is to approach the flux using an upwind finite volume scheme. This technique consists of approximating the flux by means of the values u_K, u_L , and $dV_{KL} := \tau_{KL} (v_L - v_K)$ which are available in the neighborhood of the interface σ_{KL} . To do this, we use a numerical convection flux function G of arguments $(a, b, c) \in \mathbb{R}^3$ which is required to satisfy the properties :

$$\left\{ \begin{array}{l} \text{(a) } G(\cdot, b, c) \text{ is nondecreasing for all } b, c \in \mathbb{R}, \\ \text{and } G(a, \cdot, c) \text{ is nonincreasing for all } a, c \in \mathbb{R}; \\ \text{(b) } G(a, b, c) = -G(b, a, -c) \text{ for all } a, b, c \in \mathbb{R}; \\ \text{(c) } G(a, a, c) = \Psi(a) c \text{ for all } a, c \in \mathbb{R}; \\ \text{(d) there exists } C > 0 \text{ such that} \\ \quad \forall a, b, c \in \mathbb{R} \quad |G(a, b, c)| \leq C (|a| + |b|) |c|; \\ \text{(e) there exists a modulus of continuity } \omega : \mathbb{R}^+ \rightarrow \mathbb{R}^+ \text{ such that} \\ \quad \forall a, b, a', b', c \in \mathbb{R} \quad |G(a, b, c) - G(a', b', c)| \leq |c| \omega(|a - a'| + |b - b'|). \end{array} \right. \quad (2.40)$$

Remark 2. Note that the assumptions on the dependence of G on a, b are standard (see, e.g. [31]). For instance, the assumptions (a), (b), and (c) ensure respectively the maximum principle on the discrete solutions, the conservativity of the numerical flux, and the consistency of the numerical flux on each interface. Practical examples of numerical convective flux functions can be found in [31].

In our context, one possibility to construct the numerical flux G satisfying (2.40) is to split Ψ into the nondecreasing part Ψ_\uparrow and the nonincreasing part Ψ_\downarrow :

$$\Psi_\uparrow(z) := \int_0^z (\Psi'(s))^+ ds, \quad \Psi_\downarrow(z) := - \int_0^z (\Psi'(s))^- ds.$$

Herein, $s^+ = \max(s, 0)$ and $s^- = \max(-s, 0)$. Then we take

$$G(a, b, c) = c^+ (\Psi_\uparrow(a) + \Psi_\downarrow(b)) - c^- (\Psi_\uparrow(b) + \Psi_\downarrow(a)). \quad (2.41)$$

Discretization of system (2.29)

We are now in a position to discretize problem (2.29). We denote by \mathcal{D} an admissible discretization of Q_{t_f} , which consists of an admissible mesh of Ω and a uniform time step $\Delta t > 0$. We give to the parameter h the sense of

$$\max \left\{ \Delta t, \max_{K \in \mathcal{T}} \text{diam}(K), \max_{K \in \mathcal{T}} \max_{L \in \mathcal{N}(K)} d_{KL} \right\}.$$

A finite volume scheme for the discretization of the problem (2.29) is given by the following set of equations : for all $K \in \mathcal{T}$

$$u_K^0 = \frac{1}{|K|} \int_K u_0(\mathbf{x}) d\mathbf{x}, \quad v_K^0 = \frac{1}{|K|} \int_K v_0(\mathbf{x}) d\mathbf{x}, \quad (2.42)$$

and for all $K \in \mathcal{T}$ and $n \in \llbracket 0 \cdots N \rrbracket$

$$\begin{aligned}
 |K| \frac{u_K^{n+1} - u_K^n}{\Delta t} - d_1 \sum_{L \in N(K)} \frac{|\sigma_{KL}|}{d_{KL}} (A(u_L^{n+1}) - A(u_K^{n+1})) \\
 + \sum_{L \in N(K)} G(u_K^{n+1}, u_L^{n+1}; dV_{KL}^{n+1}) = |K| \mu u_K^{n+1} \left(1 - \frac{u_K^{n+1}}{u_c} \right), \quad (2.43) \\
 |K| \frac{v_K^{n+1} - v_K^n}{\Delta t} - d_2 \sum_{L \in N(K)} \frac{|\sigma_{KL}|}{d_{KL}} (v_L^{n+1} - v_K^{n+1}) = |K| (\alpha u_K^n - \beta v_K^{n+1}),
 \end{aligned}$$

with the discrete unknowns $U = (u_K^{n+1})_{K \in \mathcal{T}}$ and $V = (v_K^{n+1})_{K \in \mathcal{T}}$, $n \in \llbracket 0 \cdots N \rrbracket$.

Notice that the discrete zero-flux boundary conditions are implicitly contained in equations (2.43).

Indeed, we have for all $K \in \mathcal{T}$, the contribution of $\partial\Omega \cap \partial T$ to the approximation of $\int_{\partial K} \nabla v \cdot \eta$,

$\int_{\partial K} \nabla A(u) \cdot \eta$ is zero, in compliance with the last equation of system (2.29).

The convergence analysis of the finite volume scheme (2.42)–(2.43) is carried out in [36] for the nondegenerate case and in [7] for the degenerate case. We refer to those references to make sure that the scheme (2.42)–(2.43) converges towards the weak solution of problem (2.29).

2.5.3 Numerical results

In this section, we show two numerical simulations of model (2.29) in a two dimensional space to illustrate our theoretical results and demonstrate the patterns generated by the system. To do this, we use the Newton algorithm to approach the solution of the nonlinear system defined by the first equation of (2.43). This algorithm is coupled with a bigradient method to solve linear systems arising from the Newton algorithm process.

We simulate the model in a two dimensional domain $\Omega = (0, 5) \times (0, 5)$ for which we consider a nonuniform admissible grid (see Figure 2.6). We consider the following data for the numerical

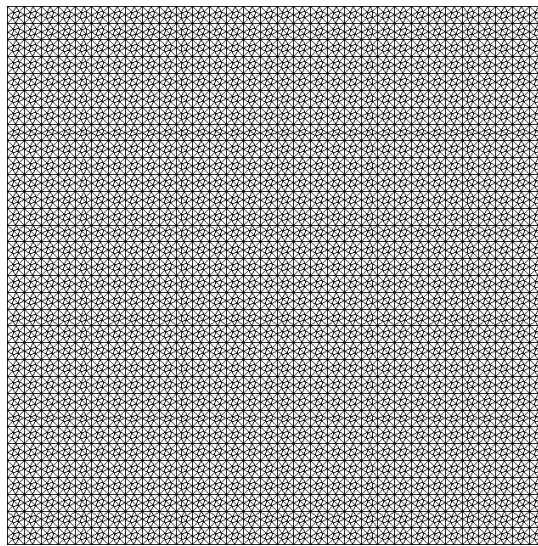


FIGURE 2.6 – Admissible mesh with 14 346 triangles.

test $d_1 = 0.1$, $d_2 = 1.0$, $u_c = 0.25$, $\bar{u} = 1.0$, $\mu = 0.5$, $\alpha = 10.0$, $\beta = 10.0$, $\chi = 20$, and $\gamma = 2$.

Further, we consider a small time step $\Delta t = 0.05$ and in the definition of the numerical fluxes G defined by (2.40), we consider

$$\Psi_{\uparrow}(z) = \Psi(\min\{z, u_m\}) \text{ and } \Psi_{\downarrow}(z) = \Psi(\max\{z, u_m\}) - \Psi(u_m)$$

where $u_m = \frac{\bar{u}}{\sqrt[\gamma]{\gamma + 1}}$.

Initial conditions and boundary conditions

Initially the cell density is setting as a small perturbation around the homogeneous steady state ; in our test, we take $u(0, \mathbf{x}, \mathbf{y}) = u_s + \varepsilon \cos(n\pi\mathbf{x}) \cos(m\pi\mathbf{y})$ where $u_s = u_c = 0.25$, $\varepsilon = 0.001$, $m = 6$, and $n = 0$. While for the chemoattractant, we consider a random perturbation around the steady state. Fig. 2.7 shows the plot of the initial condition for the cell density and for the chemoattractant ; the red regions correspond to the highest cell density so the cells are aggregated into four columns “*four stripes*”. For the boundary conditions, we impose zeros-flux boundary conditions.

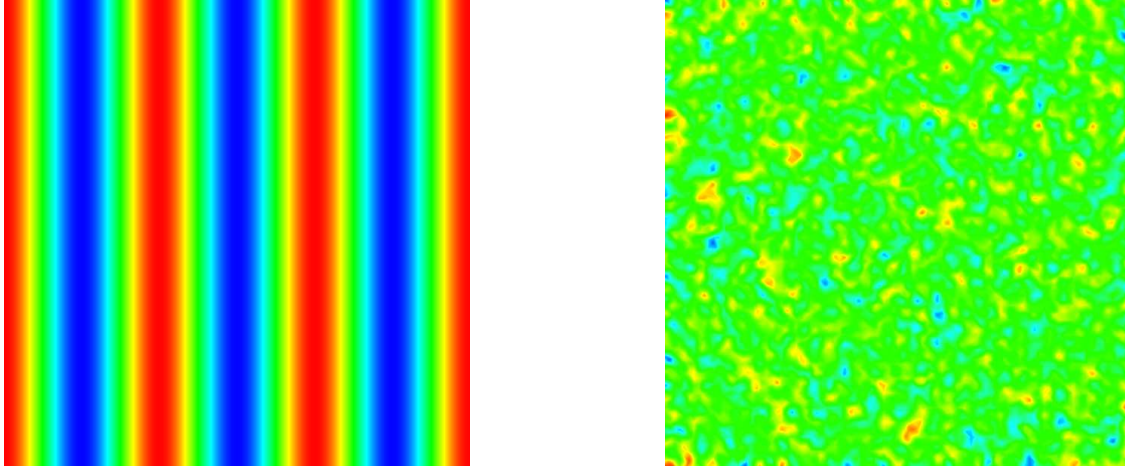


FIGURE 2.7 – Initial condition for each of the cell density $0.249 \leq u \leq 0.251$ (left) with a small perturbation around the steady state, and for the chemoattractant (right) with random perturbation around v_s .

Spatial evolution of the cell density

Figures 2.8-2.10 show the evolution of the cell density in space for different moments. We note that at $t = 0.1$ s the cell density u begins to disperse with response to the gradient of the chemoattractant v but the cell density remains distributed around the four columns. Then after half a second, the spatial patterns begin to form and the aggregations start to merge and emerge ($t = 1$ s) and then to form uniform spot patterns as at $t = 14$ s. Next, after 228s, these spot patterns stabilize and finally form 13 spatially inhomogeneous stable patterns as we can see in Fig 2.10, where the two plots show the same patterns for different moments.

Time evolution of the cell density

Here we are interested in the time evolution of the cell density at fixes points in the last plot of Fig. 2.10. Indeed, we want to show that the cell density stabilize at a certain moment ; thus we verify the theoretical results about the pattern formation using an appropriate numerical scheme.

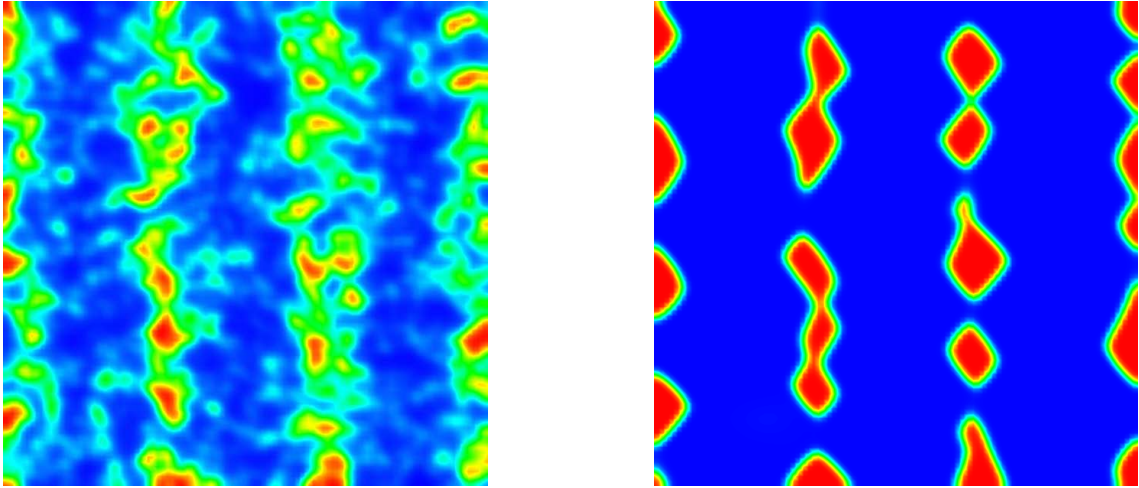


FIGURE 2.8 – Pattern formation for the full chemotaxis model (2.29) at time $t = 0.1$ s with $1.76 \times 10^{-2} \leq u \leq 0.97$ (left) and at time $t = 0.6$ s with $3.65 \times 10^{-5} \leq u \leq 0.99$ (right).

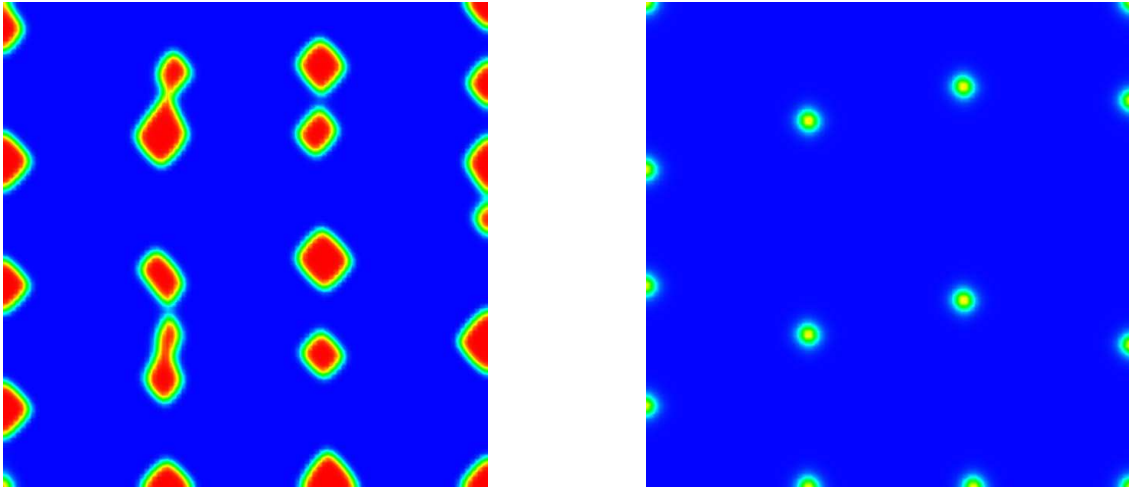


FIGURE 2.9 – Pattern formation for the full chemotaxis model (2.29) at time $t = 1$ s with $2.46 \times 10^{-6} \leq u \leq 0.99$ (left) and at time $t = 14$ s with 1.3×10^{-3} (right).

Figure 2.11 shows the evolution of the cell density at point P_1 (3.375; 2.025) in the blue line and at point P_2 (4 : 975; 2 : 825) in the red line. We observe that the cell density at the two points increases and decreases with response to the gradient of the chemoattractant which plays an essential role to stop the aggregation of the cells.

Next, the cell density stabilizes for both points when t is greater or equal to 178s, we get the same result for the others spot patterns ; we do not provide them here for convenience. Consequently, this stabilizing proves that the volume-filling chemotaxis model (2.29) generates stationary spatial patterns.

Test2 Now we consider a second numerical test for the generation of spatial patterns of the volume-filling chemotaxis model. Here, for the initial condition of the cell density, we take a random distribution of the form $u(0, \mathbf{x}, \mathbf{y}) = u_s + \varepsilon \times \text{rand}(0, 1)$ where $u_s = 0.25$, $\varepsilon = 0.01$, and rand is function which returns a pseudo-random number from a uniform distribution between 0 and 1. We take the same data as the previous test except for the domain size which replaced by $\Omega = (0, 10) \times (0, 10)$. Figures 2.12–2.14 show the evolution of the cell density in space for

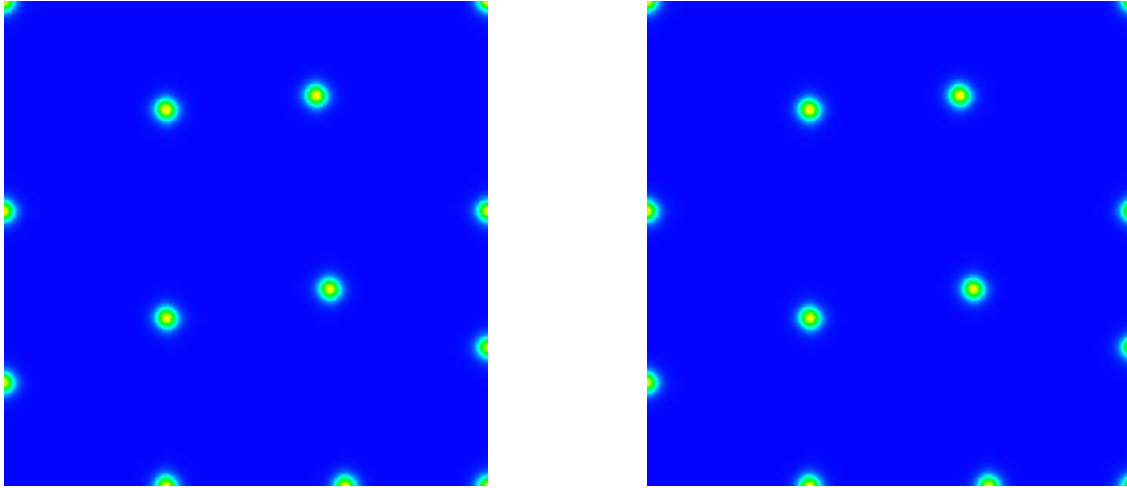


FIGURE 2.10 – Pattern formation for the full chemotaxis model (2.29) at time $t = 242$ s with $1.45 \times 10^{-3} \leq u \leq 0.79$ (left) and at time $t = 316$ s with $1.45 \times 10^{-3} \leq u \leq 0.79$ (right).

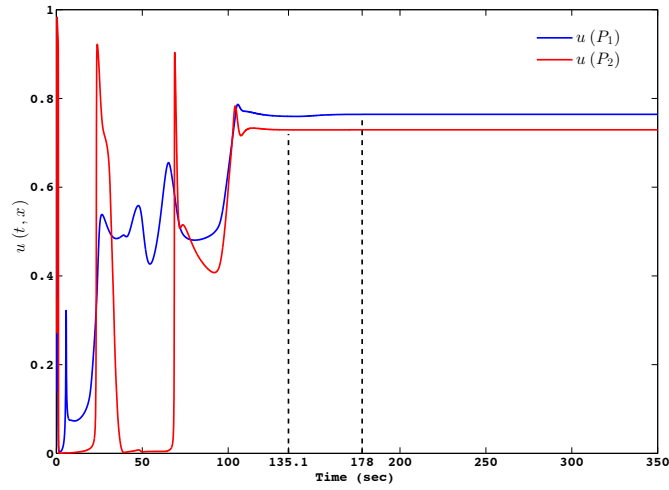


FIGURE 2.11 – Evolution of the cell density with respect to the time at two different points P_1 and P_2 .

different moments. As before, we observe the emerging process and then we get stationary spot patterns.

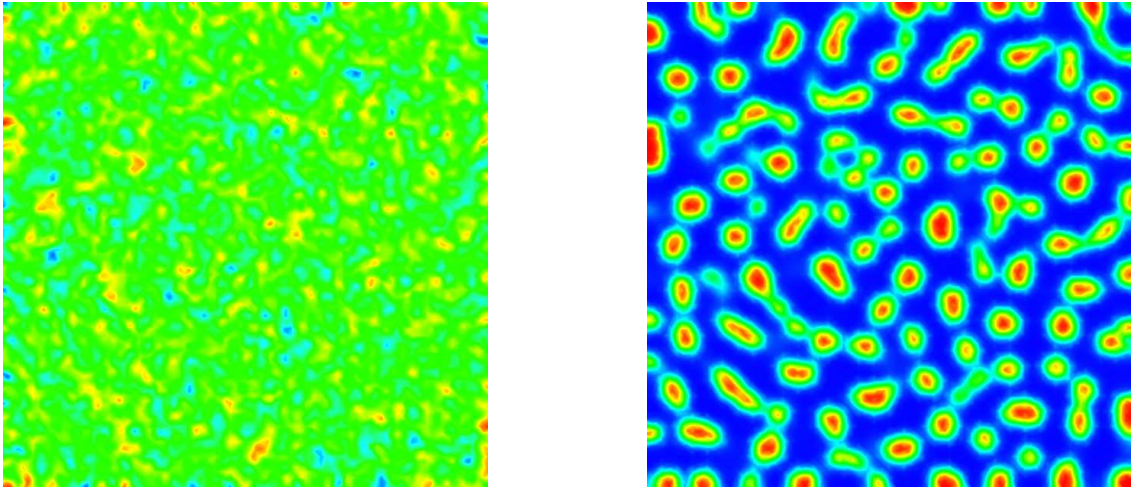


FIGURE 2.12 – Pattern formation for the full chemotaxis model (2.29) at time $t = 0$ s with $0.25 \leq u \leq 0.2599$ (left) and at time $t = 4$ s with $3.3 \times 10^{-3} \leq u \leq 0.99$ (right).

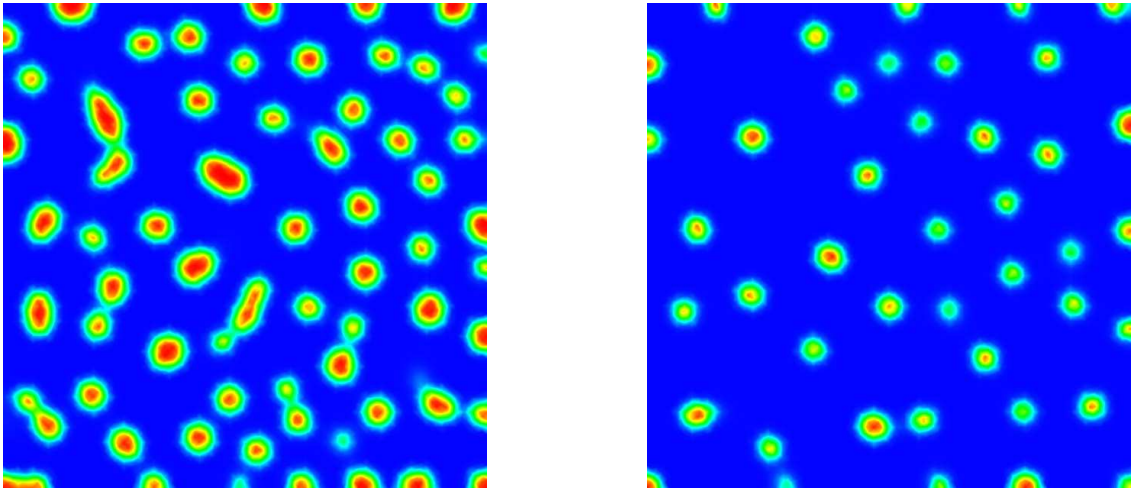


FIGURE 2.13 – Pattern formation for the full chemotaxis model (2.29) at time $t = 10$ s with $1.87 \times 10^{-5} \leq u \leq 0.99$ (left) and at time $t = 240$ s with $2.14 \times 10^{-5} \leq u \leq 0.97$ (right).

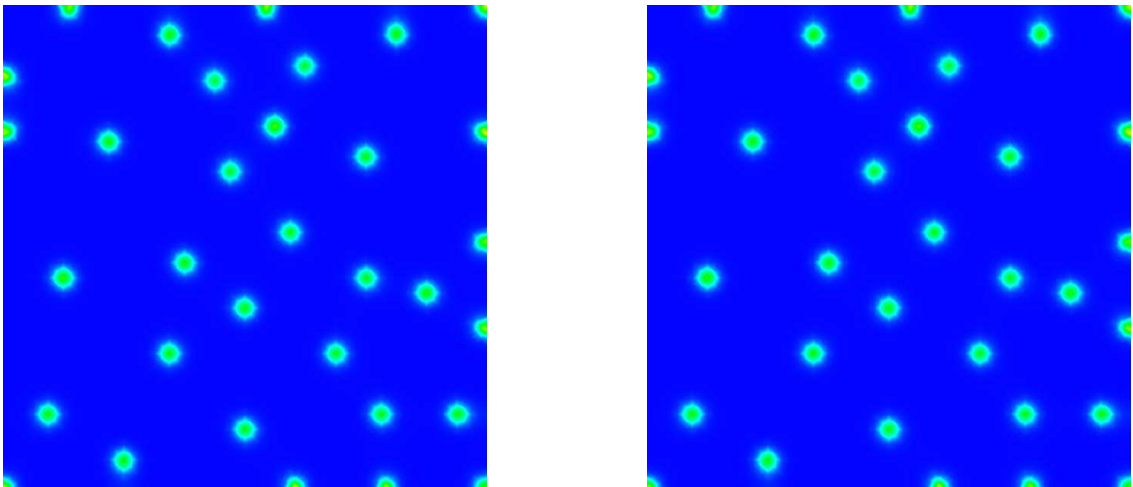


FIGURE 2.14 – Pattern formation for the full chemotaxis model (2.29) at time $t = 400$ s with $1.07 \times 10^{-3} \leq u \leq 0.691$ (left) and at time $t = 500$ s with $1.07 \times 10^{-3} \leq u \leq 0.691$ (right).

CVFE scheme for the capture of patterns for a volume-filling chemotaxis model

Sommaire

3.1	Introduction	53
3.2	Volume-filling chemotaxis model	55
3.3	CVFE discretization of the continuous problem	58
3.3.1	CVFE scheme for the modified Keller-Segel model	60
3.3.2	Main result	63
3.4	<i>A priori</i> analysis of discrete solutions	64
3.4.1	Nonnegativity of $v_{\mathcal{M},\Delta t}$, confinement of $u_{\mathcal{M},\Delta t}$	64
3.4.2	Discrete <i>a priori</i> estimates	66
3.4.3	Existence of a discrete solution	69
3.5	Compactness estimates on discrete solutions	71
3.5.1	Time translate estimate	71
3.5.2	Space translate estimate	73
3.6	Convergence of the CVFE scheme	74
3.7	Numerical simulation in two-dimensional space	79

3.1 Introduction

Patterns are the solutions of a reaction–diffusion system which are stable in time and stationary inhomogeneous in space, while pattern formation in mathematics refers to the process that, by changing a bifurcation parameter, the spatially homogeneous steady states lose stability to spatially inhomogeneous perturbations, and stable inhomogeneous solutions arise.

The pattern formation has been successfully applied to bacteria (see e.g. [77]) where we investigate specific and necessary parameters to obtain stationary distribution of the disease. Also, it

has been applied to skin pigmentation patterns [59] to understand the diversity of patterns on the animal coat pattern, and many other examples.

The pattern formation depends on two key properties : the first is to apply the seminal idea of Turing [76] for a reaction–diffusion system and consequently determine the bifurcation parameters for the generation of stationary inhomogeneous spatial patterns (also called Turing Patterns), and the second is to apply a robust scheme to numerically investigate and capture the generation of spatio-temporal patterns. One of the most popular reaction–diffusion systems that can generate spatial patterns is the chemotaxis model.

Chemotaxis is the feature movement of a cell along a chemical concentration gradient either towards the chemical stimulus, and in this case the chemical is called chemoattractant, or away from the chemical stimulus and then the chemical is called chemorepellent. The mathematical analysis of chemotaxis models shows a plenitude of spatial patterns such as the chemotaxis models applied to skin pigmentation patterns [63, 67] that lead to aggregations of one type of pigment cell into a striped spatial pattern. Other models have been successfully applied to the aggregation patterns in an epidemic disease [9], tumor growth [19], angiogenesis in tumor progression [14], and many other examples. Theoretical and mathematical modeling of chemotaxis dates to the pioneering works of Patlak in the 1950s [69] and Keller and Segel in the 1970s [51, 52]. The review article by Horstmann [47] provides a detailed introduction into the mathematics of the Keller-Segel model for chemotaxis.

In this chapter, we present and study a numerical scheme for the capture of spatial patterns for a nonlinear degenerate volume-filling chemotaxis model over a general mesh, and with inhomogeneous and anisotropic diffusion tensors. Recently, the convergence analysis of a finite volume scheme for a degenerate chemotaxis model over a homogeneous domain has been studied by Andreianov *et al.* [7], where the diffusion tensor is considered to be proportional to the identity matrix, and the mesh used for the space discretization is assumed to be admissible in the sense of satisfying the orthogonality condition as in [31]. The upwind finite volume method used for the discretization of the convective term ensures stability and is extremely robust and computationally inexpensive. However, standard finite volume scheme does not permit handling anisotropic diffusion on general meshes, even if the orthogonality condition is satisfied. The reason for this is that there is no straightforward way to apply the finite volume scheme to problems with full diffusion tensors. Various “multi-point” schemes, where the approximation of the flux through an edge involves several scalar unknowns, have been proposed, see for e.g. [33, 32] for the so-called SUCHI scheme, [25, 22] for the so-called gradient scheme, and [3] for the development of the so-called DDFV schemes. However, such schemes require using more points than the classical 4 points for triangular meshes in space dimension two, making the schemes less robust and more susceptible to numerical instabilities.

To handle the discretization of the anisotropic diffusion, it is well-known that the finite element method allows for an easy discretization of the diffusion term with a full tensor. However, it is well-established that numerical instabilities may arise in the convection-dominated case. To avoid these instabilities, the theoretical analysis of the *control volume finite element method* has been carried out for the case of degenerate parabolic problems with full diffusion tensors. Schemes with mixed conforming piecewise linear finite elements on triangles for the diffusion term and finite volumes on dual elements were proposed and studied in [15, 20, 29, 2, 17] for fluid mechanics equations, are indeed quite efficient.

Afif and Amaziane analyzed in [2] the convergence of a vertex-centered finite volume scheme for a nonlinear and degenerate convection-diffusion equation modeling a flow in porous media and without reaction term. This scheme consists of a discretization of the Laplacian by the piecewise

linear conforming finite element method (see also [64, 37]), the effectiveness of this scheme was tested in benchmarks of FVCA series of conferences [42]. Cariaga *et al.* in [17] considered the same scheme for a reaction–diffusion–convection system, where the velocity of the fluid flow is considered to be constant in the convective term.

The intention of this chapter is to extend the ideas of [7, 2, 17] to a fully nonlinear degenerate parabolic system modeling the effect of volume-filling for chemotaxis. In order to discretize this class of systems, we discretize the diffusion term by means of piecewise linear conforming finite element. The other terms are discretized by means of a finite volume scheme on a dual mesh (Donald mesh), with an upwind discretization of the numerical flux of the convective term to ensure the stability and the maximum principle of the scheme, where the dual mesh is constructed around each vertex of every triangle of the primary mesh.

The rest of this chapter is organized as follows. In Section 3.2, we introduce the chemotaxis model based on realistic biological assumptions, which incorporates the effect of volume-filling mechanism and leads to a nonlinear degenerate parabolic system. In Section 3.3, we derive the *control volume finite element scheme*, where an upwind finite volume scheme is used for the approximation of the convective term, and a standard P1-finite element method is used for the diffusive term. In Section 3.4, by assuming that the transmissibility coefficients are nonnegative, we prove the discrete maximum principle and give the *a priori* estimates on the discrete solutions. In Section 3.5, we show the compactness of the set of discrete solutions by deriving estimates on difference of time and space translates for the approximate solutions. Next, in Section 3.6, using the Kolmogorov relative compactness theorem, we prove the convergence of a sequence of the approximate solutions, and we identify the limits of the discrete solutions as weak solutions of the parabolic system proposed in Section 3.2. In the last section, we present some numerical simulations to capture the generation of spatial patterns for the volume-filling chemotaxis model with different tensors. These numerical simulations are obtained with our *control volume finite element scheme*.

3.2 Volume-filling chemotaxis model

We are interested in the *control volume finite element scheme* for a nonlinear, degenerate parabolic system formed by reaction–diffusion–convection equations. This system is complemented with homogeneous zero flux boundary conditions, which correspond to the physical behavior of the cells and the chemoattractant. The modified Keller-Segel system that we consider here, is very similar to that of Andreianov *et al.* [7], to which we have added tensors for the diffusion terms. Specifically, we consider the following system :

$$\begin{cases} \partial_t u - \operatorname{div} (\Lambda(\mathbf{x}) a(u) \nabla u - \Lambda(\mathbf{x}) \chi(u) \nabla v) = f(u) & \text{in } Q_{t_f}, \\ \partial_t v - \operatorname{div} (D(\mathbf{x}) \nabla v) = g(u, v) & \text{in } Q_{t_f}, \end{cases} \quad (3.1)$$

with the boundary conditions on $\Sigma_{t_f} := \partial\Omega \times (0, t_f)$ given by

$$(\Lambda(\mathbf{x}) a(u) \nabla u - \Lambda(\mathbf{x}) \chi(u) \nabla v) \cdot \mathbf{n} = 0, \quad D(\mathbf{x}) \nabla v \cdot \mathbf{n} = 0, \quad (3.2)$$

where, \mathbf{n} is the unit outward normal vector at the boundary $\partial\Omega$ of the domain Ω .

The initial conditions are given by :

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad v(\mathbf{x}, 0) = v_0(\mathbf{x}), \quad \mathbf{x} \in \Omega. \quad (3.3)$$

Herein, $Q_{t_f} := \Omega \times (0, t_f)$, $t_f > 0$ is a fixed time, and Ω is an open bounded polygonal domain in \mathbb{R}^2 , with Lipschitz boundary $\partial\Omega$. The initial conditions u_0 and v_0 satisfy : $u_0, v_0 \in$

$L^\infty(\Omega)$ such that $0 \leq u_0(\mathbf{x}) \leq 1$ and $v_0(\mathbf{x}) \geq 0$, for all $\mathbf{x} \in \Omega$.

In the above model, the density of the cell-population and the chemoattractant (or repellent) concentration are represented by $u = u(\mathbf{x}, t)$ and $v = v(\mathbf{x}, t)$ respectively. Next, $a(u)$ is a density-dependent diffusion coefficient, and $\Lambda(\mathbf{x})$ is the diffusion tensor in a heterogeneous medium. Furthermore, the function $\chi(u)$ is the chemoattractant sensitivity, and $D(\mathbf{x})$ is the diffusion tensor for v . The function $f(u)$ describes cell proliferation and cell death, it is usually considered to follow the logistic growth with certain carrying capacity u_c which represents the maximum density that the environment can support (e.g. see [10]). The function $g(u, v)$ describes the rates of production and degradation of the chemoattractant; here, we assume it is the linear function given by

$$g(u, v) = \alpha u - \beta v, \quad \alpha, \beta \geq 0. \quad (3.4)$$

Painter and Hillen [66] introduced the mechanistic description of the volume-filling effect. In the volume-filling effect, it is assumed that particles have a finite volume and that cells cannot move into regions that are already filled by other cells. First, we give a brief derivation of the model below, where in addition we consider the elastic cell property; that is, we consider that the cells are deformable and elastic and can squeeze into openings.

The derivation of the model begins with a master equation for a continuous-time and discrete space-random walk introduced by Othmer and Stevens [74], that is

$$\frac{\partial u_i}{\partial t} = \mathcal{C}_{i-1}^+ u_{i-1} + \mathcal{C}_{i+1}^- u_{i+1} - (\mathcal{C}_i^+ + \mathcal{C}_i^-) u_i, \quad (3.5)$$

where \mathcal{C}_i^\pm are the transitional probabilities per unit of time for one-step jump to $i \pm 1$. Herein, we shall equate the probability distribution above with the cell density.

In the volume-filling approach, and in the context of chemotaxis, the probability of a cell making a jump is assumed to depend on additional factors, such as the external concentration of the chemotactic agent and the availability of space into which the cells can squeeze and move. Therefore, we consider in the transition probability the fact that the cells can detect a local gradient as well as squeeze into openings. We take

$$\mathcal{C}_i^\pm = q(u_{i\pm 1}) (\theta + \delta [\tau(v_{i\pm 1}) - \tau(v_i)]), \quad (3.6)$$

where $q(u)$ is a nonlinear function representing the squeezing probability of a cell finding space at its neighboring location, θ and δ are constants, and τ represents the mechanism of the signal detection of the chemical concentration (for more details see [66, 74, 45]). It is assumed that only a finite number of cells, \bar{u} , can be accommodated at any site, and the function q is stipulated by the following condition :

$$q(\bar{u}) = 0 \quad \text{and} \quad 0 < q(u) \leq 1 \quad \text{for } 0 \leq u < \bar{u}.$$

Clearly, a possible choice for the squeezing probability q is a nonlinear function (see [78] for more details), defined by

$$q(u) = \begin{cases} 1 - \left(\frac{u}{\bar{u}}\right)^\gamma, & 0 \leq u \leq \bar{u}, \\ 0, & u > \bar{u}, \end{cases} \quad (3.7)$$

where $\gamma \geq 1$ denotes the squeezing exponent. The case $\gamma = 1$ corresponds to the interpretation that cells are solid blocks. However, the cells are elastic and can squeeze into openings. Thus the squeezing probability should be considered as a nonlinear function.

Substituting equation (3.6) into the master equation (3.5) and assuming that the cell density can diffuse in a heterogeneous manner in the space we get the first equation of system (3.1), with the associated coefficients

$$a(u) = d_1 (q(u) - q'(u) u), \quad \chi(u) = \zeta u q(u), \quad (3.8)$$

where d_1 and ζ are two positive constants.

In this chapter, we are interested in system (3.1) modeling the volume-filling chemotaxis process in the general case and for which we set $\bar{u} = 1$. Furthermore, we assume that the functions f and χ are continuous and satisfy :

$$f(0) = f(\bar{u}) = 0, \quad \chi(0) = \chi(\bar{u}) = 0. \quad (3.9)$$

Now, we state the main assumptions made about the system (3.1)–(3.3)

- (A1) $a : [0, 1] \rightarrow \mathbb{R}$ is a continuous function such that : $a(0) \geq 0$, $a(1) \geq 0$ and $a(u) > 0$ for all $0 < u < 1$.
- (A2) $\chi : [0, 1] \rightarrow \mathbb{R}$ is a differentiable function such that : $\chi(0) = \chi(1) = 0$ and $\chi(u) > 0$ for all $0 < u < 1$.
- (A3) The diffusion tensors Λ and D are two bounded, uniformly positive symmetric tensors on Ω , that is : $\forall \mathbf{w} \neq 0$, there exists two nonnegative constants T_- and T_+ such that $0 < T_- |\mathbf{w}|^2 \leq \langle T(\mathbf{x})\mathbf{w}, \mathbf{w} \rangle \leq T_+ |\mathbf{w}|^2 < \infty$, $T = \Lambda$ or D .
- (A4) The cell proliferation function $f : \mathbb{R} \rightarrow \mathbb{R}$ is a continuous function such that : $f(0) \geq 0$ and $f(1) \leq 0$.
- (A5) The initial conditions u_0 and v_0 are two functions lying into $L^\infty(\Omega)$ such that : $0 \leq u_0 \leq 1$ and $v_0 \geq 0$.

In the sequel, we use the nonlinear Lipschitz continuous nondecreasing function $A : [0, 1] \rightarrow \mathbb{R}$ defined by

$$A(u) = \int_0^u a(s) \, ds. \quad (3.10)$$

Definition 3.1 (weak solution). Under the assumptions (A1)–(A5). The couple of measurable functions (u, v) is said to be a weak solution of the system (3.1)–(3.3) if

$$\begin{aligned} 0 \leq u(\mathbf{x}, t) \leq 1, \quad v(\mathbf{x}, t) \geq 0 \text{ for almost everywhere } (\mathbf{x}, t) \in Q_{t_f}, \\ A(u) \in L^2(0, t_f; H^1(\Omega)), \\ v \in L^\infty(Q_{t_f}) \cap L^2(0, t_f; H^1(\Omega)), \end{aligned}$$

and for all $\varphi, \psi \in \mathcal{D}(\bar{\Omega} \times [0, t_f])$

$$\begin{aligned} - \int_{\Omega} u_0(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x} - \iint_{Q_{t_f}} u \partial_t \varphi \, d\mathbf{x} \, dt + \iint_{Q_{t_f}} \Lambda(\mathbf{x}) \nabla A(u) \cdot \nabla \varphi \, d\mathbf{x} \, dt \\ - \iint_{Q_{t_f}} \Lambda(\mathbf{x}) \chi(u) \nabla v \cdot \nabla \varphi \, d\mathbf{x} \, dt = \iint_{Q_{t_f}} f(u) \varphi(\mathbf{x}, t) \, d\mathbf{x} \, dt, \\ - \int_{\Omega} v_0(\mathbf{x}) \psi(\mathbf{x}, 0) \, d\mathbf{x} - \iint_{Q_{t_f}} v \partial_t \psi \, d\mathbf{x} \, dt \\ + \iint_{Q_{t_f}} D(\mathbf{x}) \nabla v \cdot \nabla \psi \, d\mathbf{x} \, dt = \iint_{Q_{t_f}} g(u, v) \psi(\mathbf{x}, t) \, d\mathbf{x} \, dt. \end{aligned}$$

3.3 CVFE discretization of the continuous problem

We are placed in the case where the boundary of the domain occupied by the cells and the chemoattractant is fixed over the time. We consider that $\Omega \subset \mathbb{R}^2$ is the domain occupied by the cells, and that $\Omega_h \subset \Omega$ is the approached polygonal domain of Ω . The construction of the approached solution requires the introduction of two different space discretizations of the domain Ω_h . The first discretization namely the *primal mesh* consists of simple polygons that form a partition of the domain Ω_h . The approach adopted here is oriented to the summits (“Vertex Centered approach”), the objective of this technique is to construct the solution at the summits of the *primal mesh*. To do this, a new partition of the domain is defined in such a way that every vertex of the *primal mesh* is included into only a polygon of the new partition. The mesh resulting from the second partition is called *dual mesh* and the polygons that compose it are called *dual control volumes*.

Definition 3.2. (*Primal and dual mesh*). Let Ω be an open bounded polygonal connected subset of \mathbb{R}^2 . A primal finite volume mesh of Ω is a triplet $(\mathcal{T}, \mathcal{E}, \mathcal{P})$, where \mathcal{T} is a family of disjoint open polygonal convex subsets of Ω called control volumes, \mathcal{E} is a family of subsets of $\overline{\Omega}$ contained in straight lines of \mathbb{R}^2 with strictly positive one-dimensional measure, called the edges of the control volumes, and \mathcal{P} is a family of points of Ω satisfying the following properties :

1. The closure of the union of all the control volumes is $\overline{\Omega}$, i.e. $\overline{\Omega} = \bigcup_{K \in \mathcal{T}} \overline{K}$.
2. For any $K \in \mathcal{T}$, there exists a subset \mathcal{E}_K of \mathcal{E} such that $\partial K = \overline{K} \setminus K = \bigcup_{\sigma \in \mathcal{E}_K} \overline{\sigma}$. Furthermore, $\mathcal{E} = \bigcup_{K \in \mathcal{T}} \mathcal{E}_K$.
3. There exists a Donald dual mesh $\mathcal{M} := \{M_i, i = 1, \dots, N_s\}$ associated with the triangulation $\mathcal{T} := \{K_i, i = 1, \dots, N_e\}$. For each triangle $K \in \mathcal{T}$, we connect the barycenter \mathbf{x}_K with the midpoint of each edge $\sigma \in \mathcal{E}_K$, and thus the barycenter of $K \in \mathcal{T}$ is such that $\mathbf{x}_K := \bigcap_{M \cap K \neq \emptyset} \partial M \in K$. We denote by \mathbf{x}_M the center of each dual control volume $M \in \mathcal{M}$ defined by $\mathbf{x}_M := \bigcap_{K \cap M \neq \emptyset} \partial K \in M$. For each interface of the dual control volume M , we denote by $\sigma_{M,M'}^K := \partial M \cap \partial M' \cap K$ the line segment between the points \mathbf{x}_K and the midpoint of the line segment $[\mathbf{x}_M, \mathbf{x}_{M'}]$ and let $\mathcal{L} := \{\sigma \in \partial M \setminus \partial \Omega, M \in \mathcal{M}\}$ be the set of all interior sides. Finally we denote by \mathcal{M}^{int} and \mathcal{M}^{ext} the set of all interior and all boundary dual control volumes respectively. We refer to Fig. 3.1 for an illustration of the primal triangular mesh \mathcal{T} and its corresponding Donald dual mesh \mathcal{M} .

In the sequel, we use the following notations. For any $M \in \mathcal{M}$, $|M|$ is the area of M . The set of neighbors of M is denoted by $\mathcal{N}(M) := \{M' \in \mathcal{M} / \exists \sigma \in \mathcal{L}, \overline{\sigma} = \overline{M} \cap \overline{M'}\}$, and we designate by $d_{M,M'}$ the distance between the centers of M and M' . We define the mesh size by $h := \text{size}(\mathcal{M}) = \sup_{M \in \mathcal{M}} \text{diam}(M)$ and make the following shape regularity assumption on the family of triangulations $\{\mathcal{T}_h\}_h$:

$$\text{There exists a positive constant } \kappa_{\mathcal{T}} \text{ such that : } \min_{K \in \mathcal{T}_h} \frac{|K|}{\text{diam}(K)^2} \geq \kappa_{\mathcal{T}}, \quad \forall h > 0. \quad (3.11)$$

For the time discretization, we do not impose any restriction on the time step, for that we consider a uniform time step $\Delta t \in (0, t_f)$. We take $N \in \mathbb{N}^*$ such that $N := \max\{n \in \mathbb{N} / n \Delta t < t_f\}$, and we denote $t_n = n \Delta t$, for $n \in \{0, \dots, N+1\}$, so that $t_0 = 0$, and $t_{N+1} = t_f$.

We define the following finite-dimensional spaces :

$$\begin{aligned} \mathcal{H}_{\mathcal{T}} &:= \{\varphi \in C^0(\overline{\Omega}) ; \varphi|_K \in \mathbb{P}_1, \forall K \in \mathcal{T}\} \subset H^1(\Omega), \\ \mathcal{H}_{\mathcal{T}}^0 &:= \{\varphi \in \mathcal{H}_{\mathcal{T}} ; \varphi(\mathbf{x}_M) = 0, \quad \forall M \in \mathcal{M}^{\text{ext}}\}. \end{aligned}$$

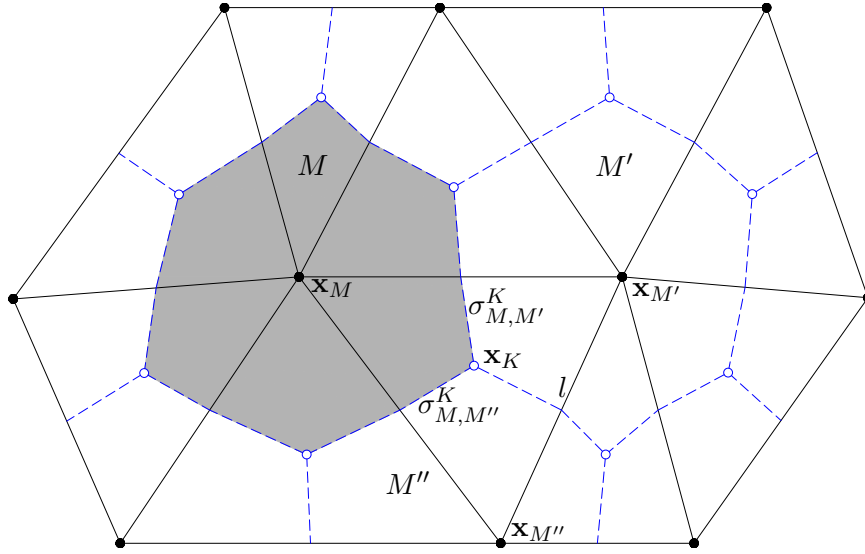


FIGURE 3.1 – Primal triangular mesh \mathcal{T} and Donald dual mesh \mathcal{M} : control volumes, centers, interfaces.

The canonical basis of $\mathcal{H}_{\mathcal{T}}$ is spanned by the shape functions $(\varphi_M)_{M \in \mathcal{M}}$, such that $\varphi_M(\mathbf{x}_{M'}) = \delta_{M,M'}$ for all $M' \in \mathcal{M}$, δ being the Kronecker delta. The approximations in these spaces are conforming since $\mathcal{H}_{\mathcal{T}} \subset H^1(\Omega)$. We equip $\mathcal{H}_{\mathcal{T}}$ with the semi-norm

$$\|u_{\mathcal{T}}\|_{\mathcal{H}_{\mathcal{T}}}^2 := \int_{\Omega} |\nabla u_{\mathcal{T}}|^2 d\mathbf{x}, \quad \forall u_{\mathcal{T}} \in \mathcal{H}_{\mathcal{T}},$$

which becomes a norm on $\mathcal{H}_{\mathcal{T}}^0$.

The classical finite elements \mathbb{P}_1 associated to the vertex \mathbf{x}_{M_i} ($i = 1, \dots, N_s$), where N_s is the total number of vertices, is defined by

$$\varphi_{M_i}(\mathbf{x}_{M_j}) = \delta_{ij}, \text{ where } \varphi_{M_i} \text{ is continuous and piecewise } \mathbb{P}_1 \text{ per triangle.}$$

Let w_M^n be an expected approximation of $w(\mathbf{x}_M, t_n)$, where $w \equiv u$ or v . Thus, the discrete unknowns are denoted by $\{w_M^n / M \in \mathcal{M}, n \in \{0, \dots, N+1\}\}$.

Definition 3.3. (Discrete functions). Using the values of (u_M^{n+1}, v_M^{n+1}) , $\forall M \in \mathcal{M}$ and $n \in \{0, \dots, N\}$, we determine two approximate solutions by means of the *control volume finite element scheme* :

- (i) A finite volume solution $(u_{\mathcal{M},\Delta t}, v_{\mathcal{M},\Delta t})$ defined as piecewise constant on the dual control volumes in space and piecewise constant in time, such that :

$$\begin{aligned} (u_{\mathcal{M},\Delta t}(\mathbf{x}, 0), v_{\mathcal{M},\Delta t}(\mathbf{x}, 0)) &= (u_M^0, v_M^0) & \forall \mathbf{x} \in \overset{\circ}{M}, M \in \mathcal{M}, \\ (u_{\mathcal{M},\Delta t}(\mathbf{x}, t), v_{\mathcal{M},\Delta t}(\mathbf{x}, t)) &= (u_M^{n+1}, v_M^{n+1}) & \forall \mathbf{x} \in \overset{\circ}{M}, M \in \mathcal{M}, \forall t \in (t_n, t_{n+1}], \end{aligned}$$

where u_M^0 (resp. v_M^0) represents the mean value of the function u_0 (resp. v_0). The discrete space of these functions namely *discrete control volumes space* is denoted by $\mathcal{X}_{\mathcal{M},\Delta t}$.

- (ii) A finite element solution $v_{\mathcal{T},\Delta t}$ as a function continuous and piecewise \mathbb{P}_1 per triangle in space and piecewise constant in time, such that :

$$\begin{aligned} v_{\mathcal{T},\Delta t}(\mathbf{x}, 0) &= v_{\mathcal{T}}^0(\mathbf{x}) & \forall \mathbf{x} \in \Omega, \\ v_{\mathcal{T},\Delta t}(\mathbf{x}, t) &= v_{\mathcal{T}}^{n+1}(\mathbf{x}) & \forall \mathbf{x} \in \Omega, \forall t \in (t_n, t_{n+1}], \end{aligned}$$

where $v_{\mathcal{T}}^{n+1}(\mathbf{x}) := \sum_{M \in \mathcal{M}} v_M^{n+1} \varphi_M(\mathbf{x})$ and $v_{\mathcal{T}}^0(\mathbf{x}) := \sum_{M \in \mathcal{M}} v_M^0 \varphi_M(\mathbf{x})$. The discrete space of these functions namely *discrete finite elements space* is denoted by $\mathcal{H}_{\mathcal{T}, \Delta t}$.

The function A is nonlinear, we denote by $A_{\mathcal{T}, \Delta t} = A_{\mathcal{T}}(u_{\mathcal{T}, \Delta t})$ the corresponding finite element reconstruction in $\mathcal{H}_{\mathcal{T}, \Delta t}$, and by $A_{\mathcal{M}, \Delta t} = A(u_{\mathcal{M}, \Delta t})$ the corresponding finite volume reconstruction in $\mathcal{X}_{\mathcal{M}, \Delta t}$. Specifically, we have

$$\begin{aligned} A_{\mathcal{T}}(u_{\mathcal{T}, \Delta t}(\mathbf{x}, t)) &= \sum_{M \in \mathcal{M}} A(u_M^{n+1}) \varphi_M(\mathbf{x}), \quad \forall \mathbf{x} \in \Omega, \forall t \in (t_n, t_{n+1}], \\ A(u_{\mathcal{M}, \Delta t}(\mathbf{x}, t)) &= A(u_M^{n+1}), \quad \forall \mathbf{x} \in \overset{\circ}{M}, M \in \mathcal{M}, \forall t \in (t_n, t_{n+1}]. \end{aligned}$$

3.3.1 CVFE scheme for the modified Keller-Segel model

In order to define a discretization for system (3.1), we integrate the equations of system (3.1) over the set $M \times [t_n, t_{n+1}]$ with $M \in \mathcal{M}$, then we use the Green–Gauss formula as well as the implicit order one discretization in time, we get

$$\begin{aligned} \int_M (u(\mathbf{x}, t_{n+1}) - u(\mathbf{x}, t_n)) d\mathbf{x} - \sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \Lambda \nabla A(u) \cdot \eta_{M, \sigma} d\sigma(\mathbf{x}) dt \\ + \sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \chi(u) \Lambda \nabla v \cdot \eta_{M, \sigma} d\sigma(\mathbf{x}) dt = \int_{t_n}^{t_{n+1}} \int_M f(u) d\mathbf{x} dt, \\ \int_M (v(\mathbf{x}, t_{n+1}) - v(\mathbf{x}, t_n)) d\mathbf{x} - \sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} D \nabla v \cdot \eta_{M, \sigma} d\sigma(\mathbf{x}) dt = \int_{t_n}^{t_{n+1}} \int_M g(u, v) d\mathbf{x} dt, \end{aligned} \quad (3.12)$$

where $\eta_{M, \sigma}$ is the unit normal vector outward to $\sigma \subset \partial M$, and $d\sigma(\mathbf{x})$ is the Lebesgue measure on the edge σ .

We consider now an implicit Euler scheme in time, and thus the time evolution in the first equation of system (3.12) is approximated as

$$\begin{aligned} \int_M (u(\mathbf{x}, t_{n+1}) - u(\mathbf{x}, t_n)) d\mathbf{x} &\approx \int_M (u_{\mathcal{M}, \Delta t}(\mathbf{x}, t_{n+1}) - u_{\mathcal{M}, \Delta t}(\mathbf{x}, t_n)) d\mathbf{x} \\ &= |M| (u_M^{n+1} - u_M^n). \end{aligned}$$

Note that $f(u)$ is a nonlinear function. We denote by $f(u_{\mathcal{M}, \Delta t})$ the corresponding piecewise constant reconstruction in $\mathcal{X}_{\mathcal{M}, \Delta t}$, then the reaction term is approximated as

$$\int_{t_n}^{t_{n+1}} \int_M f(u(\mathbf{x}, t)) d\mathbf{x} dt \approx \int_{t_n}^{t_{n+1}} \int_M f(u_{\mathcal{M}, \Delta t}(\mathbf{x}, t)) d\mathbf{x} dt = |M| \Delta t f(u_M^{n+1}).$$

On the other hand, we distinguish two kinds of approximation in space. The first consists of considering the finite element approach to handle the diffusion term, and the second consists of using a classical upstream finite volume approach.

Let us focus on the discretization of the diffusion term of the first equation of system (3.12), we have

$$\sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \Lambda \nabla A(u) \cdot \eta_{M, \sigma} d\sigma(\mathbf{x}) dt \approx \Delta t \sum_{\sigma \subset \partial M} \int_{\sigma} \Lambda \nabla A_{\mathcal{T}}(u_{\mathcal{T}, \Delta t}(\mathbf{x}, t_{n+1})) \cdot \eta_{M, \sigma} d\sigma(\mathbf{x}). \quad (3.13)$$

The diffusion tensor $\Lambda(\mathbf{x})$ is taken constant per triangle, we denote by Λ_K the mean value of the function $\Lambda(\mathbf{x})$ over the triangle $K \in \mathcal{T}$, then one rewrites the right hand side of equation (3.13) as

$$\begin{aligned} \Delta t \sum_{K, K \cap M \neq \emptyset} \sum_{\sigma \subset \partial M \cap K} \Lambda_K \nabla A_{\mathcal{T}}(u_{\mathcal{T}, \Delta t}(\mathbf{x}, t_{n+1}))|_K \cdot \eta_{M, \sigma} |\sigma| \\ = \Delta t \sum_{K, K \cap M \neq \emptyset} \frac{1}{2} \nabla A_{\mathcal{T}}(u_{\mathcal{T}, \Delta t}(\mathbf{x}, t_{n+1}))|_K \cdot {}^t \Lambda_K \eta_{K, l} |l| \end{aligned} \quad (3.14)$$

where $l \in \mathcal{E}_K$ such that $M \cap l = \emptyset$, and $\eta_{K, l}$ denotes the unit normal vector outward to the edge l . For the transition between the first and the second line in approximation (3.14), we have used a geometric property (see Lemma A.3), that is

$$\sum_{\sigma \subset \partial M \cap K} \eta_{M, \sigma} |\sigma| = \frac{1}{2} \eta_{K, l} |l|, \quad \forall K \in \mathcal{T} \text{ such that } K \cap M \neq \emptyset.$$

According to the definition of the approximate function $\nabla A_{\mathcal{T}}(u_{\mathcal{T}, \Delta t})$, one has

$$\nabla A_{\mathcal{T}}(u_{\mathcal{T}, \Delta t}(\mathbf{x}, t_{n+1}))|_K = \nabla \left(\sum_{M \in \mathcal{M}} A(u_M^{n+1}) \varphi_M(\mathbf{x}) \right)|_K = \sum_{M \in \mathcal{M}} A(u_M^{n+1}) \nabla \varphi_M(\mathbf{x})|_K. \quad (3.15)$$

Note that, the \mathbb{P}_1 -finite element bases are expressed in barycentric coordinates, thus

$$\sum_{M', M' \cap K \neq \emptyset} \varphi_{M'}(\mathbf{x})|_K = 1, \text{ and } \sum_{M', M' \cap K \neq \emptyset} \nabla \varphi_{M'}(\mathbf{x})|_K = 0,$$

consequently, we have

$$\nabla A_{\mathcal{T}}(u_{\mathcal{T}, \Delta t}(\mathbf{x}, t_{n+1}))|_K = \sum_{M', M' \cap K \neq \emptyset} (A(u_{M'}^{n+1}) - A(u_M^{n+1})) \nabla \varphi_{M'}|_K. \quad (3.16)$$

We note that, for a given $K \in \mathcal{T}$, we have (see Appendix A)

$$\nabla \varphi_M|_K = \frac{-|l|}{2|K|} \eta_{K, l}, \quad \forall M \in \mathcal{M} \text{ such that } M \cap K \neq \emptyset. \quad (3.17)$$

Let us now introduce the transmissibility coefficient between M and M' defined, for all $K \in \mathcal{T}$ such that $K \cap M \neq \emptyset$ and $K \cap M' \neq \emptyset$, by

$$\Lambda_{M, M'}^K = - \int_K \Lambda(\mathbf{x}) \nabla \varphi_M(\mathbf{x}) \cdot \nabla \varphi_{M'}(\mathbf{x}) d\mathbf{x}. \quad (3.18)$$

As a consequence of equations (3.17)–(3.18), one has

$$\sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \Lambda \nabla A(u) \cdot \eta_{M, \sigma} d\sigma(\mathbf{x}) dt \approx \Delta t \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} \Lambda_{M, M'}^K (A(u_{M'}^{n+1}) - A(u_M^{n+1})).$$

Next, we have to approximate the convective term in the first equation of system (3.12). For that, we consider a classical upstream finite volume scheme according to the normal component of the gradient of the chemoattractant v on the interfaces. So,

$$\begin{aligned} \sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \chi(u) \Lambda \nabla v \cdot \eta_{M, \sigma} d\sigma(\mathbf{x}) dt \\ \approx \Delta t \sum_{K, K \cap M \neq \emptyset} \sum_{\sigma \subset \partial M \cap K} \int_{\sigma} \Lambda(\mathbf{x}) \chi(u_{\mathcal{M}, \Delta t}(\mathbf{x}, t_{n+1})) \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t_{n+1}) \cdot \eta_{M, \sigma} d\sigma(\mathbf{x}). \end{aligned}$$

In order to approximate the convective flux on each interface, let us firstly introduce an example of approximation in the case where the function χ is nondecreasing. For that, we consider the interface $\sigma_{M,M'}^K$, and write

$$\int_{\sigma_{M,M'}^K} \chi(u_{M,\Delta t}(\mathbf{x}, t_{n+1})) \Lambda(\mathbf{x}) \nabla v_{T,\Delta t}(\mathbf{x}, t_{n+1}) \cdot \eta_\sigma d\sigma(\mathbf{x}) \\ \approx \begin{cases} |\sigma_{M,M'}^K| \chi(u_M^{n+1}) dV_{M,M'}^K, & \text{if } dV_{M,M'}^K \geq 0, \\ |\sigma_{M,M'}^K| \chi(u_{M'}^{n+1}) dV_{M,M'}^K, & \text{if } dV_{M,M'}^K \leq 0, \end{cases}$$

where $dV_{M,M'}^K$ represents an approximation of the gradient of v on the interface $\sigma_{M,M'}^K$, defined by

$$dV_{M,M'}^K = \sum_{M'', M'' \cap K \neq \emptyset} v_{M''}^{n+1} \nabla \varphi_{M''|_K} \cdot {}^t \Lambda_K \eta_{M,M'}^K. \quad (3.19)$$

Thus, the convective term is approximated as

$$\sum_{\sigma \subset \partial M} \int_{t_n}^{t_{n+1}} \int_{\sigma} \chi(u) \Lambda \nabla v \cdot \eta_{M,\sigma} d\sigma(\mathbf{x}) dt \\ \approx \begin{cases} \Delta t \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} \chi(u_M^{n+1}) \Lambda_{M,M'}^K (v_{M'}^{n+1} - v_M^{n+1}), & \text{if } dV_{M,M'}^K \geq 0, \\ \Delta t \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} \chi(u_{M'}^{n+1}) \Lambda_{M,M'}^K (v_{M'}^{n+1} - v_M^{n+1}), & \text{if } dV_{M,M'}^K \leq 0, \end{cases}$$

where $\Lambda_{M,M'}^K$ represents the transmissibility coefficient between M and M' , given by equation (3.18).

In the general case, we use numerical convection flux functions G of arguments $(a, b, c) \in \mathbb{R}^3$ which are required to satisfy the following properties :

$$\left\{ \begin{array}{l} \text{(a) } G(\cdot, b, c) \text{ is nondecreasing for all } b, c \in \mathbb{R}, \\ \text{and } G(a, \cdot, c) \text{ is nonincreasing for all } a, c \in \mathbb{R}; \\ \text{(b) } G(a, b, c) = -G(b, a, -c) \text{ for all } a, b, c \in \mathbb{R}; \\ \text{(c) } G(a, a, c) = \chi(a) c \text{ for all } a, c \in \mathbb{R}; \\ \text{(d) there exists } C > 0 \text{ such that} \\ \quad \forall a, b, c \in \mathbb{R} \quad |G(a, b, c)| \leq C(|a| + |b|)|c|; \\ \text{(e) there exists a modulus of continuity } \omega : \mathbb{R}^+ \rightarrow \mathbb{R}^+ \text{ such that} \\ \quad \forall a, b, a', b', c \in \mathbb{R} \quad |G(a, b, c) - G(a', b', c)| \leq |c| \omega(|a - a'| + |b - b'|). \end{array} \right. \quad (3.20)$$

In our context, one possibility to construct a numerical flux G satisfying conditions (3.20) is to split χ into the nondecreasing part χ_\uparrow and the nonincreasing part χ_\downarrow :

$$\chi_\uparrow(z) := \int_0^z (\chi'(s))^+ ds, \quad \chi_\downarrow(z) := - \int_0^z (\chi'(s))^- ds.$$

Herein, $s^+ = \max(s, 0)$ and $s^- = \max(-s, 0)$. Then we take

$$G(a, b, c) = c^+ (\chi_\uparrow(a) + \chi_\downarrow(b)) - c^- (\chi_\uparrow(b) + \chi_\downarrow(a)). \quad (3.21)$$

Notice that in the case χ has a unique local (and global) maximum at the point $\tilde{u} \in (0, 1)$, we have

$$\chi_\uparrow(z) = \chi(\min\{z, \tilde{u}\}) \quad \text{and} \quad \chi_\downarrow(z) = \chi(\max\{z, \tilde{u}\}) - \chi(\tilde{u}).$$

For the discretization of the second equation of system (3.12), we define the transmissibility coefficient $D_{M,M'}^K$ by

$$D_{M,M'}^K = \int_K D(\mathbf{x}) \nabla \varphi_M(\mathbf{x}) \cdot \nabla \varphi_{M'}(\mathbf{x}) d\mathbf{x}. \quad (3.22)$$

then we follow the same lines as for the discretization of the first equation.

We are now in a position to discretize problem (3.1)–(3.3). We denote by \mathcal{D} a discretization of Q_{t_f} , which consists of a *primal finite element triangulation* \mathcal{T} of Ω and its corresponding *Donald dual mesh* \mathcal{M} and a time step $\Delta t > 0$.

A *control volume finite element scheme* for the discretization of problem (3.1)–(3.3) is given by the following set of equations : for all $M \in \mathcal{M}$,

$$u_M^0 = \frac{1}{|M|} \int_M u_0(\mathbf{x}) d\mathbf{x}, \quad v_M^0 = \frac{1}{|M|} \int_M v_0(\mathbf{x}) d\mathbf{x}, \quad (3.23)$$

and for all $M \in \mathcal{M}$ and all $n \in \{0, \dots, N\}$,

$$\begin{aligned} |M| \frac{u_M^{n+1} - u_M^n}{\Delta t} - \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} \Lambda_{M,M'}^K (A(u_{M'}^{n+1}) - A(u_M^{n+1})) \\ + \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} |\sigma_{M,M'}^K| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) = |M| f(u_M^{n+1}), \end{aligned} \quad (3.24)$$

$$|M| \frac{v_M^{n+1} - v_M^n}{\Delta t} - \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} D_{M,M'}^K (v_{M'}^{n+1} - v_M^{n+1}) = |M| g(u_M^n, v_M^{n+1}). \quad (3.25)$$

We recall that the unknowns are $U = (u_M^{n+1})_{M \in \mathcal{M}}$ and $V = (v_M^{n+1})_{M \in \mathcal{M}}$, $n \in \{0, \dots, N\}$, and that $dV_{M,M'}^K$ is defined in equation (3.19), and the transmissibility coefficients $\Lambda_{M,M'}^K$ and $D_{M,M'}^K$ are given by equation (3.18) and equation (3.22) respectively. Notice that the discrete zero-flux boundary conditions are implicitly contained in equations (3.24)–(3.25). The contribution of $\partial\Omega \cap \partial M$ to the approximation of $\int_{\partial M} D \nabla v \cdot \eta d\sigma(\mathbf{x})$ and $\int_{\partial M} \Lambda (\nabla A(u) - \chi(u) \nabla v) \cdot \eta d\sigma(\mathbf{x})$ is zero, in compliance with equation (3.2).

In this chapter, we assume that all the transmissibility coefficients $\Lambda_{M,M'}^K$ and $D_{M,M'}^K$ are nonnegative :

$$\Lambda_{M,M'}^K \geq 0 \text{ and } D_{M,M'}^K \geq 0, \quad \forall M, M' \in \mathcal{M}, \forall K \in \mathcal{T}. \quad (3.26)$$

For the case where Λ and D are isotropic, and if all the triangles of \mathcal{T} have only acute angles, the condition (3.26) holds. This was used for example in [18] to prove the convergence of the finite element approximation towards the renormalized solution of an elliptic equation.

3.3.2 Main result

Let $(\mathcal{T}_m)_{m \geq 1}$ be a sequence of triangulations of Ω such that

$$h_m = \max_{K \in \mathcal{T}_m} \text{diam}(K) \rightarrow 0 \text{ as } m \rightarrow \infty.$$

We assume that there exists a constant $\kappa > 0$ such that

$$\kappa \mathcal{T}_m \leq \kappa, \quad \forall m \geq 1.$$

As before, a sequence of dual meshes $(\mathcal{M}_m)_{m \geq 1}$ is given.

Let $(N_m)_m$ be an increasing sequence of integers, then we define the corresponding sequence of time steps $(\Delta t_m)_m$ such that $\Delta t_m \rightarrow 0$ as $m \rightarrow \infty$. The intention of this chapter is to prove the following main result.

Theorem 3.4. *Let $(u_{\mathcal{M}_m, \Delta t_m}, v_{\mathcal{M}_m, \Delta t_m})_m$ be a sequence of solutions to the scheme (3.23)–(3.25), such that $0 \leq u_{\mathcal{M}_m, \Delta t_m} \leq 1$ and $0 \leq v_{\mathcal{M}_m, \Delta t_m}$ for almost everywhere in Q_{t_f} , then*

$$u_{\mathcal{M}_m, \Delta t_m} \rightarrow u \text{ and } v_{\mathcal{M}_m, \Delta t_m} \rightarrow v \quad \text{a.e. in } Q_{t_f} \text{ as } m \rightarrow \infty,$$

where the couple (u, v) is a weak solution to the system (3.1)–(3.3) in the sense of Definition 3.1.

3.4 A priori analysis of discrete solutions

In this section, we prove the discrete maximum principle, then we establish the *a priori* estimates necessary to prove the existence of a solution to the discrete problem (3.23)–(3.25) and the convergence of the scheme towards the weak solution.

In the sequel, we denote by C a “generic” constant, which need not have the same value through the proofs.

3.4.1 Nonnegativity of $v_{\mathcal{M}, \Delta t}$, confinement of $u_{\mathcal{M}, \Delta t}$

We aim to prove the following lemma which is a basis to the analysis that we are going to perform.

Lemma 3.5. *Let $(u_M^n, v_M^n)_{M \in \mathcal{M}, n \in \{0, \dots, N+1\}}$ be a solution of the CVFE scheme (3.23)–(3.25). Under the nonnegativity of transmissibility coefficients assumption (3.26), we have for all $M \in \mathcal{M}$, and all $n \in \{0, \dots, N+1\}$, $0 \leq u_M^n \leq 1$ and $0 \leq v_M^n$. Moreover, there exists a positive constant $\rho = \|v_0\|_\infty + \alpha t_f$, such that $v_M^n \leq \rho$, for all $n \in \{0, \dots, N+1\}$.*

Proof. Let us show by induction on n that for all $M \in \mathcal{M}$, $u_M^n \geq 0$. The claim is true for $n = 0$. We argue by induction that for all $M \in \mathcal{M}$, the claim is true up to order n . Consider a dual control volume M such that $u_M^{n+1} = \min \{u_{M'}^{n+1}\}_{M' \in \mathcal{M}}$, we want to show that $u_M^{n+1} \geq 0$. We consider equation (3.24) corresponding to the aforementioned dual control volume M , reorganize the summation over the edges and multiply it by $-(u_M^{n+1})^-$ where for all real r , $r = r^+ - r^-$ with $r^+ = \max(r, 0)$ and $r^- = \max(-r, 0)$. This yields

$$\begin{aligned} & -|M| \frac{u_M^{n+1} - u_M^n}{\Delta t} (u_M^{n+1})^- + \sum_{\sigma_{M, M'}^K \subset \partial M} \Lambda_{M, M'}^K (A(u_{M'}^{n+1}) - A(u_M^{n+1})) (u_M^{n+1})^- \\ & - \sum_{\sigma_{M, M'}^K \subset \partial M} |\sigma_{M, M'}^K| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M, M'}^K) (u_M^{n+1})^- = -f(u_M^{n+1}) (u_M^{n+1})^-. \end{aligned} \quad (3.27)$$

Here, we use the extension by $f(0) \geq 0$ of the continuous function f for $u \leq 0$, then the right hand side of equation (3.27) is less or equal to zero.

The function A is nondecreasing then $A(u_{M'}^{n+1}) - A(u_M^{n+1}) \geq 0$ since $u_M^{n+1} \leq u_{M'}^{n+1}$. Furthermore, the assumption $\Lambda_{M, M'}^K \geq 0$ implies that

$$\sum_{\sigma_{M, M'}^K \subset \partial M} \Lambda_{M, M'}^K (A(u_{M'}^{n+1}) - A(u_M^{n+1})) (u_M^{n+1})^- \geq 0.$$

From the assumptions on the numerical flux G , the function G is nonincreasing with respect to the second variable $u_{M'}^{n+1}$, then using the extension by zero of the continuous function χ (recall that $\chi(u) = 0$ for $u \leq 0$), we get

$$\begin{aligned} G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M, M'}^K) (u_M^{n+1})^- & \leq G(u_M^{n+1}, u_M^{n+1}; dV_{M, M'}^K) (u_M^{n+1})^- \\ & = dV_{M, M'}^K \chi(u_M^{n+1}) (u_M^{n+1})^- = 0. \end{aligned}$$

Now, using the identity $u_M^{n+1} = (u_M^{n+1})^+ - (u_M^{n+1})^-$ and the nonnegativity of u_M^n , we deduce from equation (3.27) that $(u_M^{n+1})^- = 0$. According to the choice of the dual control volume M , then $\min \{u_{M'}^{n+1}\}_{M' \in \mathcal{M}}$ is nonnegative; this implies that

$$0 \leq u_M^{n+1} \leq u_{M'}^{n+1} \text{ for all } n \in \{0, \dots, N\} \text{ and all } M' \in \mathcal{M}.$$

In order to prove that $u_M^n \leq 1$ for all $n \in \{0, \dots, N+1\}$ and all $M \in \mathcal{M}$, we argue by induction that for all $M \in \mathcal{M}$, the claim $u_M^n \leq 1$ is true. We take the dual control volume M such that $u_M^{n+1} = \max (u_{M'}^{n+1})_{M' \in \mathcal{M}}$, we want to show that $u_M^{n+1} \leq 1$.

For the mentioned claim, we multiply equation (3.24) by $(u_M^{n+1} - 1)^+$, one gets

$$\begin{aligned} |M| \frac{u_M^{n+1} - u_M^n}{\Delta t} (u_M^{n+1} - 1)^+ - \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (A(u_M^{n+1}) - A(u_{M'}^{n+1})) (u_M^{n+1} - 1)^+ \\ + \sum_{\sigma_{M,M'}^K \subset \partial M} |\sigma_{M,M'}^K| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) (u_M^{n+1} - 1)^+ = f(u_M^{n+1}) (u_M^{n+1} - 1)^+. \end{aligned} \quad (3.28)$$

Since A is nondecreasing function, and $u_{M'}^{n+1} \leq u_M^{n+1}$, we get $A(u_{M'}^{n+1}) - A(u_M^{n+1}) \leq 0$. This implies that

$$- \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (A(u_M^{n+1}) - A(u_{M'}^{n+1})) (u_M^{n+1} - 1)^+ \geq 0. \quad (3.29)$$

Next, we use again the fact that the numerical flux function G is nonincreasing with respect to the second variable u_M^{n+1} and consistence (see (b) and (c) in properties (3.20)) to deduce, using the extension by zero of the continuous function χ for $u \geq 1$, that

$$\begin{aligned} G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) (u_M^{n+1} - 1)^+ &\geq G(u_M^{n+1}, u_M^{n+1}; dV_{M,M'}^K) (u_M^{n+1} - 1)^+ \\ &= dV_{M,M'}^K \chi(u_M^{n+1}) (u_M^{n+1} - 1)^+ = 0. \end{aligned} \quad (3.30)$$

Now, we rely on the extension by $f(1) \leq 0$ on $[1, +\infty[$ of the production term source f in the right hand side of (3.28), thus $f(u_M^{n+1}) (u_M^{n+1} - 1)^+ \leq 0$.

Using the above estimates (3.29)–(3.30) into equation (3.28), one has

$$\begin{aligned} |M| \frac{u_M^{n+1} - u_M^n}{\Delta t} (u_M^{n+1} - 1)^+ &= \frac{|M|}{\Delta t} \left((u_M^{n+1} - 1) (u_M^{n+1} - 1)^+ \right. \\ &\quad \left. - (u_M^n - 1) (u_M^{n+1} - 1)^+ \right) \leq 0. \end{aligned} \quad (3.31)$$

Using again the identity $(u_M^{n+1} - 1) = (u_M^{n+1} - 1)^+ - (u_M^{n+1} - 1)^-$, and that $u_M^n \leq 1$, one can deduce from estimate (3.31) that $(u_M^{n+1} - 1)^+ = 0$. Consequently, we obtain

$$u_{M'}^{n+1} \leq u_M^{n+1} \leq 1 \text{ for all } n \in \{0, \dots, N\} \text{ and all } M' \in \mathcal{M}.$$

The proof of nonnegativity of v_M^n , $M \in \mathcal{M}$, $n \in \{0, \dots, N+1\}$, follows the same lines as in the proof for the nonnegativity of u_M^n , since $-g(u_M^n, v_M^{n+1})(u_M^{n+1})^- = -\alpha |M| u_M^n (v_M^{n+1})^- + \beta |M| v_M^{n+1} (v_M^{n+1})^- \leq 0$. For simplicity, we do not provide it here.

Let us now focus on the last claim concerning the existence of a constant ρ such that $v_M^n \leq \rho$ for all $M \in \mathcal{M}$ and all $n \in \{0, \dots, N+1\}$. We set $\rho_n := \|v_0\|_\infty + n\alpha\Delta t$, and suppose that

$v_M^n \leq \rho_n$, $\forall M \in \mathcal{M}$ (the claim holds for $n = 0$). We want to show that $v_M^{n+1} \leq \rho_{n+1}$, for that we take the dual control volume M such that $v_M^{n+1} = \max \{v_{M'}^{n+1}\}_{M' \in \mathcal{M}}$. Using scheme (3.25) and the fact that $\rho_{n+1} = \rho_n + \alpha \Delta t$, one has

$$\begin{aligned} |M| \frac{v_M^{n+1} - \rho_{n+1}}{\Delta t} + |M| \beta v_M^{n+1} - \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} D_{M, M'}^K (v_{M'}^{n+1} - v_M^{n+1}) \\ = \alpha |M| (u_M^n - 1) + |M| \frac{v_M^n - \rho_n}{\Delta t}. \end{aligned} \quad (3.32)$$

Multiplying equation (3.32) by $(v_M^{n+1} - \rho_{n+1})^+$, one can deduce that $v_M^{n+1} \leq \rho_{n+1} \leq \rho$, for all $n \in \{0, \dots, N\}$. This ends the proof of the lemma. \square

3.4.2 Discrete *a priori* estimates

We derive in the next proposition, the main uniform estimates on the discrete gradient of the cell density $A(u)$ and the discrete gradient of the chemical concentration v . These estimates are necessary to obtain later results on the compactness as well as the existence of discrete solutions to the discrete problem

Proposition 3.6. *Let $(u_M^{n+1}, v_M^{n+1})_{M \in \mathcal{M}, n \in \{0, \dots, N\}}$, be a solution of the control volume finite element scheme (3.23)–(3.25). Under assumption (3.11) and assumption (3.26), there exists a constant $C > 0$, depending only on Ω , t_f , $\|v_0\|_\infty$, $\kappa_{\mathcal{T}}$, α , f , and on the constant in (3.20)(d) such that*

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M, M'}^K \subset \partial M} \Lambda_{M, M'}^K |A(u_M^{n+1}) - A(u_{M'}^{n+1})|^2 \\ + \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M, M'}^K \subset \partial M} D_{M, M'}^K |v_M^{n+1} - v_{M'}^{n+1}|^2 \leq C, \end{aligned} \quad (3.33)$$

consequently, for all $A_{\mathcal{T}}(u_{\mathcal{T}}^{n+1}) = \sum_{M \in \mathcal{M}} A(u_M^{n+1}) \varphi_M \in \mathcal{H}_{\mathcal{T}}$, and all $v_{\mathcal{T}}^{n+1} = \sum_{M \in \mathcal{M}} v_M^{n+1} \varphi_M \in \mathcal{H}_{\mathcal{T}}$,

$$\sum_{n=0}^N \Delta t \|v_{\mathcal{T}}^{n+1}\|_{\mathcal{H}_{\mathcal{T}}}^2 \leq C, \quad (3.34)$$

and

$$\sum_{n=0}^N \Delta t \|A_{\mathcal{T}}(u_{\mathcal{T}}^{n+1})\|_{\mathcal{H}_{\mathcal{T}}}^2 \leq C. \quad (3.35)$$

Proof. To prove estimate (3.33), we multiply equation (3.24) (resp. equation (3.25)) by $A(u_M^{n+1})$ (resp. by v_M^{n+1}), and perform a sum over $M \in \mathcal{M}$ and $n \in \{0, \dots, N\}$. This yields

$$E_{1,1} + E_{1,2} + E_{1,3} = E_{1,4} \quad \text{and} \quad E_{2,1} + E_{2,2} = E_{2,3},$$

where

$$\begin{aligned}
E_{1,1} &= \sum_{n=0}^N \sum_{M \in \mathcal{M}} |M| (u_M^{n+1} - u_M^n) A(u_M^{n+1}), \quad E_{2,1} = \sum_{n=0}^N \sum_{M \in \mathcal{M}} |M| (v_M^{n+1} - v_M^n) v_M^{n+1}, \\
E_{1,2} &= \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (A(u_M^{n+1}) - A(u_{M'}^{n+1})) A(u_M^{n+1}), \\
E_{1,3} &= \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} |\sigma_{M,M'}^K| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) A(u_M^{n+1}), \\
E_{1,4} &= \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} |M| f(u_M^{n+1}) A(u_M^{n+1}), \\
E_{2,2} &= \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} D_{M,M'}^K (v_M^{n+1} - v_{M'}^{n+1}) v_M^{n+1}, \\
E_{2,3} &= \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} |M| (\alpha u_M^n - \beta v_M^{n+1}) v_M^{n+1}.
\end{aligned}$$

We consider the Lipschitz continuous nondecreasing function $\mathcal{B} : [0, 1] \rightarrow \mathbb{R}$ defined, for every $s \in \mathbb{R}$, by $\mathcal{B}(s) := \int_0^s A(r) \, dr$; we have $\mathcal{B}''(s) = a(s) \geq 0$, so that \mathcal{B} is convex.

From the convexity of \mathcal{B} and of the function $s \rightarrow \frac{1}{2}s^2$, we have the following inequalities

$$\forall a, b \in \mathbb{R}, \quad (a - b) A(a) \geq \mathcal{B}(a) - \mathcal{B}(b), \quad (a - b) a \geq \frac{1}{2} (a^2 - b^2).$$

Using these inequalities for the terms $E_{1,1}$ and $E_{2,1}$, we obtain

$$\begin{aligned}
E_{1,1} &= \sum_{n=0}^N \sum_{M \in \mathcal{M}} |M| (u_M^{n+1} - u_M^n) A(u_M^{n+1}) \\
&\geq \sum_{n=0}^N \sum_{M \in \mathcal{M}} |M| (\mathcal{B}(u_M^{n+1}) - \mathcal{B}(u_M^n)) = \sum_{M \in \mathcal{M}} |M| (\mathcal{B}(u_M^{N+1}) - \mathcal{B}(u_M^0)), \\
E_{2,1} &= \sum_{n=0}^N \sum_{M \in \mathcal{M}} |M| (v_M^{n+1} - v_M^n) v_M^{n+1} \geq \frac{1}{2} \sum_{M \in \mathcal{M}} |M| ((v_M^{n+1})^2 - (v_M^n)^2).
\end{aligned}$$

Next, for the diffusive term $E_{1,2}$, we reorganize the sum over edges (discrete integration by parts). Then, we have

$$\begin{aligned}
E_{1,2} &= \sum_{n=0}^N \Delta t \sum_{\sigma_{M,M'}^K \in \mathcal{L}} \Lambda_{M,M'}^K (A(u_M^{n+1}) - A(u_{M'}^{n+1}))^2 \\
&= \frac{1}{2} \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (A(u_M^{n+1}) - A(u_{M'}^{n+1}))^2.
\end{aligned}$$

We estimate the convective term $E_{1,3}$, and also gather by edges, one gets

$$\begin{aligned} E_{1,3} &= \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} |\sigma_{M,M'}^K| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) A(u_M^{n+1}) \\ &= \sum_{n=0}^N \Delta t \sum_{\sigma_{M,M'}^K \in \mathcal{L}} |\sigma_{M,M'}^K| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) (A(u_M^{n+1}) - A(u_{M'}^{n+1})). \end{aligned}$$

Using the definition of the function G , the assumption (3.20)(d), the boundedness of u_M^{n+1} , and applying the weighted Young inequality, one has

$$\begin{aligned} |E_{1,3}| &\leq \sum_{n=0}^N \Delta t \sum_{\sigma_{M,M'}^K \in \mathcal{L}} |\sigma_{M,M'}^K| |G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) (A(u_M^{n+1}) - A(u_{M'}^{n+1}))| \\ &\leq E_{1,3}^1 + E_{1,3}^2, \end{aligned}$$

where

$$\begin{aligned} E_{1,3}^1 &= \frac{1}{4} \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K |A(u_M^{n+1}) - A(u_{M'}^{n+1})|^2, \\ E_{1,3}^2 &= C \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} |G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) \sigma_{M,M'}^K|^2 \\ &\leq C \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} |dV_{M,M'}^K|^2 |\sigma_{M,M'}^K|^2. \end{aligned}$$

On the other hand, using the definition (3.19) of $dV_{M,M'}^K$ and the upper bound of the diffusion tensor Λ , one gets

$$|dV_{M,M'}^K| |\sigma_{M,M'}^K| \leq C \sum_{M'', M'' \cap K \neq \emptyset} |v_{M''}^{n+1} - v_M^{n+1}| \left| \nabla \varphi_{M''|K} \right| |\sigma_{M,M'}^K|.$$

Thanks to the shape regularity assumption (3.11), one can deduce that $\left| \nabla \varphi_{M''|K} \right| |\sigma_{M,M'}^K| \leq C$. As a consequence,

$$\begin{aligned} |E_{1,3}| &\leq \frac{1}{4} \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K |A(u_M^{n+1}) - A(u_{M'}^{n+1})|^2 \\ &\quad + C \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} D_{M,M'}^K |v_{M'}^{n+1} - v_M^{n+1}|^2. \end{aligned}$$

The last estimation for the reactive term is given using definition (3.4) of g and the boundedness

of u_M^{n+1} , v_M^{n+1} , and f . Then

$$E_{2,3} = \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} |M| \left(\alpha u_M^n v_M^{n+1} - \beta (v_M^{n+1})^2 \right) \leq \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \alpha |M| v_M^{n+1} \leq \alpha \rho t_f |\Omega|.$$

$$E_{1,4} = \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} |M| f(u_M^{n+1}) A(u_M^{n+1}) \leq C t_f |\Omega|.$$

Collecting the previous inequalities, one can deduce that there exists a constant $C > 0$, independent of h and Δt , such that

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K |A(u_M^{n+1}) - A(u_{M'}^{n+1})|^2 \\ + \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} D_{M,M'}^K |v_M^{n+1} - v_{M'}^{n+1}|^2 \leq C. \end{aligned}$$

Let us focus on estimate (3.35), we denote by D_K , $K \in \mathcal{T}$ the mean value of the function D over the triangle K , then using the previous estimates as well as the assumptions on the diffusion tensor D , one has

$$\begin{aligned} C &\geq \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{\sigma_{M,M'}^K \subset \partial M} D_{M,M'}^K |v_M^{n+1} - v_{M'}^{n+1}|^2 \\ &= 2 \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} v_M^{n+1} \sum_{\sigma_{M,M'}^K \subset \partial M} D_{M,M'}^K (v_M^{n+1} - v_{M'}^{n+1}) \\ &= 2 \sum_{n=0}^N \Delta t \sum_{M \in \mathcal{M}} \sum_{K \in \mathcal{T}} |K| v_M^{n+1} \nabla \varphi_{M|_K} \cdot {}^t D_K \sum_{M' \in \mathcal{M}} v_{M'}^{n+1} \nabla \varphi_{M'|_K} \\ &= \sum_{n=0}^N \Delta t \sum_{K \in \mathcal{T}} |K| D_K \left(\sum_{M \in \mathcal{M}} v_M^{n+1} \nabla \varphi_{M|_K} \right) \cdot \left(\sum_{M' \in \mathcal{M}} v_{M'}^{n+1} \nabla \varphi_{M'|_K} \right) \\ &= \sum_{n=0}^N \Delta t \sum_{K \in \mathcal{T}} \int_K D(\mathbf{x}) \nabla v_{\mathcal{T}}^{n+1} \cdot \nabla v_{\mathcal{T}}^{n+1} d\mathbf{x} \geq 2D_- \sum_{n=0}^N \Delta t \|v_{\mathcal{T}}^{n+1}\|_{\mathcal{H}_{\mathcal{T}}}^2. \end{aligned}$$

In the same manner and using estimate (3.33), one can establish estimate (3.34). This ends the proof of the **Proposition 3.6**. \square

3.4.3 Existence of a discrete solution

The existence of a solution to the *control volume finite element scheme* will be obtained with the help of the following lemma proved in [55, 28, 75].

Lemma 3.7 (Application of Brouwer fixed-point theorem). *Let \mathcal{A} be a finite dimensional Hilbert space with inner product denoted by $[\cdot, \cdot]_{\mathcal{A}}$ and associated norm $\|\cdot\|_{\mathcal{A}}$. Let \mathcal{P} be a continuous mapping from \mathcal{A} into itself satisfying the following property : there exists $r > 0$ such that*

$$[\mathcal{P}(\xi), \xi]_{\mathcal{A}} > 0 \text{ for all } \xi \in \mathcal{A} \text{ with } \|\xi\|_{\mathcal{A}} = r.$$

Then, there exists $\xi \in \mathcal{A}$ with $\|\xi\|_{\mathcal{A}} \leq r$ such that

$$\mathcal{P}(\xi) = 0.$$

The existence of a discrete solution for the *nonlinear control volume finite element scheme* is given in the following proposition.

Proposition 3.8. *Under the shape regularity assumption (3.11) and the nonnegativity assumption (3.26), there exists at least one solution $(u_M^{n+1}, v_M^{n+1})_{(M,n) \in \mathcal{M} \times \llbracket 0 \dots N \rrbracket}$ to the discrete problem (3.23)–(3.25).*

Proof. Denote by $U_h^n = (u_M^n)_{M \in \mathcal{M}}$ and $V_h^n = (v_M^n)_{M \in \mathcal{M}}$. We will show the existence of a discrete solution by induction on n . Assume that (u_h^n, v_h^n) exists and show the existence of (u_h^{n+1}, v_h^{n+1}) . Note that the discrete equation (3.25) is a standard time-implicit finite volume discretization of a uniformly parabolic equation, where the contribution of u in the right-hand side is discretized in an explicit way. Hence, The discrete equation (3.25) can be written as a finite dimensional linear system $\mathbf{A} V_h^{n+1} = \mathbf{B}$ with respect to the unknown V_h^{n+1} . The coefficients $(\mathbf{A}_{i,j})_{i,j \in \{1, \dots, N_s\}}$ and $(\mathbf{B}_i)_{i \in \{1, \dots, N_s\}}$ of the matrix \mathbf{A} are given, for all $i \neq j \in \{1, \dots, N_s\}$, by

$$\mathbf{A}_{i,i} = 1 + \beta \Delta t + \frac{\Delta t}{|M|} \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} D_{M,M'}^K,$$

$$\mathbf{A}_{i,j} = -\frac{\Delta t}{|M|} \sum_{K, M \cap K \neq \emptyset} \sum_{M', M' \cap K \neq \emptyset} D_{M,M'}^K,$$

$$\mathbf{B}_i = \alpha \Delta t u_M^n + v_M^n.$$

It is clear that the resulting matrix \mathbf{A} of this system is symmetric and strictly diagonally dominant with nonnegative diagonal coefficients. As a consequence, \mathbf{A} is invertible and definite positive. Thus the equation (3.25) admits a unique solution V_h^{n+1} .

Let us now prove the existence of a solution to the scheme (3.24). The nonlinear function A defined by equation (3.10) is nondecreasing, hence A is invertible. We can rewrite the scheme (3.24) in terms of ω_h^i with $U_h^i = A^{-1}(\omega_h^i)$, $i \in \{0, \dots, N+1\}$. Assume that ω_h^n and V_h^{n+1} exist. We choose the component-wise inner product $[\cdot, \cdot]$ as the scalar product on $\mathbb{R}^{\text{Card}(\mathcal{M})}$. We define the mapping \mathbf{F} , that associates to the vector $\mathcal{W} = (W_M^{n+1})_{M \in \mathcal{M}}$, the following expression

$$\begin{aligned} \mathbf{F}(\mathcal{W}) = & \left(|M| \frac{A^{-1}(W_M^{n+1}) - A^{-1}(W_M^n)}{\Delta t} - \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (W_{M'}^{n+1} - W_M^{n+1}) \right. \\ & + \sum_{\sigma_{M,M'}^K \subset \partial M} |\sigma_{M,M'}^K| G(A^{-1}(W_M^{n+1}), A^{-1}(W_{M'}^{n+1}); dV_{M,M'}^K) \\ & \left. - |M| f \circ A^{-1}(W_M^{n+1}) \right)_{M \in \mathcal{M}}, \end{aligned}$$

given by equation (3.24). Now, using Lemma 3.5 and estimates (3.33), one can deduce that

$$[\mathbf{F}(\mathcal{W}), \mathcal{W}] \geq C |\mathcal{W}|^2 - C' |\mathcal{W}| - C'', \quad \text{for } |\mathcal{W}| \text{ large enough,}$$

where C, C', C'' are three nonnegative constants.

This implies, reasoning by contradiction and using Brouwer theorem (see Lemma 3.7), that, there exists $\mathcal{W} \in \mathbb{R}^{\text{Card}(\mathcal{M})}$ such that

$$\mathbf{F}(\mathcal{W}) = 0.$$

Therefore, we obtain the existence of at least one solution to the scheme (3.24). \square

3.5 Compactness estimates on discrete solutions

In this section, we derive estimates on differences of time and space translates necessary to prove the relative compactness property of the sequence of approximate solutions using Kolmogorov's theorem. Under the shape regularity assumption (3.11) and the nonnegativity of the transmissibility coefficients assumption (3.26), we give the time and space translate estimates for $A(u_{\mathcal{M}_h, \Delta t_h})$ and $v_{\mathcal{M}_h, \Delta t_h}$ given by **Definition 3.3**.

3.5.1 Time translate estimate

Lemma 3.9. *Under assumption (3.11) and assumption (3.26), there exists a positive constant $C > 0$ depending on Ω , t_f , α , u_0 and v_0 (neither on h nor on τ) such that*

$$\iint_{\Omega \times (0, t_f - \tau)} |w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt \leq C(\tau + \Delta t_h), \quad \text{for all } \tau \in (0, t_f) \quad (3.36)$$

with $w_{\mathcal{M}_h, \Delta t_h} = A(u_{\mathcal{M}_h, \Delta t_h})$, or $v_{\mathcal{M}_h, \Delta t_h}$.

Proof. Let $\tau \in (0, t_f)$ and $t \in (0, t_f - \tau)$. We consider the quantity $\Upsilon(t)$ defined by

$$\Upsilon(t) := \int_{\Omega} |A(u_{\mathcal{M}_h, \Delta t_h})(\mathbf{x}, t + \tau) - A(u_{\mathcal{M}_h, \Delta t_h})(\mathbf{x}, t)|^2 d\mathbf{x}.$$

Set $n_0(t) = [t/\Delta t_h]$ and $n_1(t) = [(t + \tau)/\Delta t_h]$, where $[x] = n$ for $x \in [n, n + 1)$, $n \in \mathbb{N}$.

The function A is nondecreasing; then using the mean value theorem, one gets the following inequality

$$\iint_{\Omega \times (0, t_f - \tau)} |A(u_{\mathcal{M}_h, \Delta t_h})(\mathbf{x}, t + \tau) - A(u_{\mathcal{M}_h, \Delta t_h})(\mathbf{x}, t)|^2 d\mathbf{x} dt \leq C \int_0^{t_f - \tau} \Upsilon(t) dt,$$

where, for almost every $t \in (0, t_f - \tau)$,

$$\Upsilon(t) = \sum_{M \in \mathcal{M}_h} |M| \left(A(u_M^{n_1(t)}) - A(u_M^{n_0(t)}) \right) \left(u_M^{n_1(t)} - u_M^{n_0(t)} \right).$$

Note that the function $\Upsilon(t)$ may be written as

$$\Upsilon(t) = \sum_{M \in \mathcal{M}_h} \left(A(u_M^{n_1(t)}) - A(u_M^{n_0(t)}) \right) \sum_{n=0}^{N-1} \chi_n(t, t + \tau) |M| (u_M^{n+1} - u_M^n), \quad (3.37)$$

where, χ_n is the characteristic function defined by

$$\chi_n(t, t + \tau) = \begin{cases} 1 & \text{if } (n + 1) \Delta t_h \in (t, t + \tau], \\ 0 & \text{if } (n + 1) \Delta t_h \notin (t, t + \tau]. \end{cases}$$

In equation (3.37), the order of the summation between n and M is changed and the scheme (3.24) is used. Hence,

$$\begin{aligned} \Upsilon(t) &= \Delta t_h \sum_{n=0}^{N-1} \chi_n(t, t + \tau) \sum_{M \in \mathcal{M}_h} \left(A(u_M^{n_1(t)}) - A(u_M^{n_0(t)}) \right) \times \\ &\quad \sum_{\sigma_{M, M'}^K \subset \partial M} \left(\Lambda_{M, M'}^K (A(u_{M'}^{n+1}) - A(u_M^{n+1})) - |\sigma_{M, M'}^K| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M, M'}^K) \right) \\ &\quad + \Delta t_h \sum_{n=0}^{N-1} \chi_n(t, t + \tau) \sum_{M \in \mathcal{M}_h} \left(A(u_M^{n_1(t)}) - A(u_M^{n_0(t)}) \right) \times |M| f(u_M^{n+1}). \end{aligned}$$

We write $\Upsilon(t) = \Delta t_h \sum_{n=0}^{N-1} \chi_n(t, t + \tau) (\Upsilon_1(t) + \Upsilon_2(t) + \Upsilon_3(t))$, where

$$\begin{aligned}\Upsilon_1(t) &:= \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K \left(A(u_M^{n_1(t)}) - A(u_M^{n_0(t)}) \right) (A(u_{M'}^{n+1}) - A(u_M^{n+1})), \\ \Upsilon_2(t) &:= \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} |\sigma_{M,M'}^K| \left(A(u_M^{n_0(t)}) - A(u_M^{n_1(t)}) \right) G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K), \\ \Upsilon_3(t) &:= \sum_{M \in \mathcal{M}_h} |M| \left(A(u_M^{n_1(t)}) - A(u_M^{n_0(t)}) \right) f(u_M^{n+1}).\end{aligned}$$

It is easy to see that

$$\sum_{n=0}^{N-1} \Delta t_h \int_0^{t_f - \tau} \chi_n(t, t + \tau) \Upsilon_3(t) dt \leq C(\tau + \Delta t_h).$$

For the first term, note that gathering by edges, using the basic triangle inequality, one has

$$\begin{aligned}\Upsilon_1(t) &\leq \frac{1}{2} \left(\sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (A(u_{M'}^{n+1}) - A(u_M^{n+1}))^2 \right. \\ &\quad + \frac{1}{2} \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K \left| A(u_M^{n_1(t)}) - A(u_{M'}^{n_1(t)}) \right|^2 \\ &\quad \left. + \frac{1}{2} \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K \left| A(u_M^{n_0(t)}) - A(u_{M'}^{n_0(t)}) \right|^2 \right).\end{aligned}$$

Using the estimates (3.33), this implies that there exists a constant $C > 0$ independent of τ and h , such that $\sum_{n=0}^{N-1} \Delta t_h \int_0^{t_f - \tau} \chi_n(t, t + \tau) \Upsilon_1(t) dt \leq C(\tau + \Delta t_h)$. Finally, applying the previous arguments, gathering by edges, and using each of the definition of G and the assumptions on it, we get

$$\begin{aligned}\Upsilon_2(t) dt &\leq \frac{C}{2} \left(\sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} \left(\left| A(u_M^{n_1(t)}) - A(u_{M'}^{n_1(t)}) \right|^2 + |v_M^{n+1} - v_{M'}^{n+1}|^2 \right) \right. \\ &\quad \left. + \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} \left(\left| A(u_M^{n_0(t)}) - A(u_{M'}^{n_0(t)}) \right|^2 + |v_M^{n+1} - v_{M'}^{n+1}|^2 \right) \right).\end{aligned}$$

We use estimate (3.33) to deduce that

$$\sum_{n=0}^{N-1} \Delta t_h \int_0^{t_f - \tau} \chi_n(t, t + \tau) \Upsilon_2(t) dt \leq C(\tau + \Delta t_h).$$

Consequently, we obtain

$$\begin{aligned} \int_0^{t_f-\tau} \Upsilon(t) dt &\leq \sum_{n=0}^{N-1} \Delta t_h \int_0^{t_f-\tau} \chi_n(t, t+\tau) (\Upsilon_1(t) + \Upsilon_3(t)) dt \\ &\quad + \sum_{n=0}^{N-1} \Delta t_h \int_0^{t_f-\tau} \chi_n(t, t+\tau) \Upsilon_2(t) dt \leq C(\tau + \Delta t_h), \end{aligned}$$

for some constant C independent of τ and h . The proof of (3.36) for $w_{\mathcal{M}_h, \Delta t_h} = v_{\mathcal{M}_h, \Delta t_h}$ follows in a similar manner. This concludes the proof of the lemma. \square

We now extend by zero the functions $w_{\mathcal{M}_h, \Delta t_h}$ and $w_{\mathcal{T}_h, \Delta t_h}$ outside of Q_{t_f} and give the time translate estimate over \mathbb{R}^3 for $w_{\mathcal{M}_h, \Delta t_h}$. Indeed, there exists a constant $C > 0$ independent of h and τ such that

$$\int_{\mathbb{R}} \int_{\mathbb{R}^2} |w_{\mathcal{M}_h, \Delta t_h}(t + \tau, \mathbf{x}) - w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt \leq C(\tau + \Delta t_h), \quad \text{for all } \tau \in (0, t_f).$$

Proof. Using the extension by zero of $w_{\mathcal{M}_h, \Delta t_h}$ outside of Q_{t_f} , one has

$$\begin{aligned} &\int_{\mathbb{R}} \int_{\mathbb{R}^2} |w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt \\ &= \int \int_{Q_{t_f-\tau}} |w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt + \int_{t_f-\tau}^{t_f} \int_{\Omega} |w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt. \end{aligned}$$

One can deduce the proof using estimate (3.36) and the L^∞ bound of the function $w_{\mathcal{M}_h, \Delta t_h}$. \square

We give now the space translate estimate on the approximate solutions.

3.5.2 Space translate estimate

Lemma 3.10. *Under assumptions (3.11) and (3.26), there exists a positive constant $C > 0$ depending on Ω , t_f , α , u_0 and v_0 (neither on h nor on τ) such that*

$$\int_0^{t_f} \int_{\mathbb{R}^2} |w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x} + \mathbf{y}, t) - w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)| d\mathbf{x} dt \leq C(|\mathbf{y}| + h), \quad \text{for all } \mathbf{y} \in \mathbb{R}^2 \quad (3.38)$$

with $w_{\mathcal{M}_h, \Delta t_h} = A(u_{\mathcal{M}_h, \Delta t_h})$, or $v_{\mathcal{M}_h, \Delta t_h}$.

Proof. We follow the same proof used in [16]. We first extend by zero the function $w_{\mathcal{T}_h, \Delta t_h}$ outside of Q_{t_f} , and then we prove that

$$\int_0^{t_f} \int_{\mathbb{R}^2} |w_{\mathcal{T}_h, \Delta t_h}(\mathbf{x} + \mathbf{y}, t) - w_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t)| d\mathbf{x} dt \leq C|\mathbf{y}|, \quad \text{for all } \mathbf{y} \in \mathbb{R}^2. \quad (3.39)$$

Thanks to Proposition 3.6, and since Q_{t_f} is of finite measure, then Hölder inequality yields

$$\|\nabla w_{\mathcal{T}_h, \Delta t_h}\|_{(L^1(Q_{t_f}))^2} \leq |Q_{t_f}|^{\frac{1}{2}} \|\nabla w_{\mathcal{T}_h, \Delta t_h}\|_{(L^2(Q_{t_f}))^2}^{\frac{1}{2}} \leq C. \quad (3.40)$$

Since A is uniformly bounded thanks to the boundedness of $u_{\mathcal{T}_h, \Delta t_h}$ and the uniform continuity of A , its extension to the whole \mathbb{R}^3 , still denoted by $w_{\mathcal{T}_h, \Delta t_h}$, belongs to $L^\infty \cap \text{BV}(\mathbb{R}^3)$ and satisfies

$$\text{TV}(w_{\mathcal{T}_h, \Delta t_h}) \leq \|\nabla w_{\mathcal{T}_h, \Delta t_h}\|_{(L^1(Q_{t_f}))^2} + \|w\|_\infty (t_f |\partial\Omega| + 2|\Omega|) < \infty, \quad (3.41)$$

where $|\partial\Omega|$ represents the length of $\partial\Omega$. The estimate (3.39) is a classical consequence of inequality (3.41) (see e.g. [13, 12]).

Now, using the extension by zero of $w_{\mathcal{M}_h, \Delta t_h}$ outside of Q_{t_f} as well as the triangle inequality, we get

$$\int_0^{t_f} \int_{\mathbb{R}^2} |w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x} + \mathbf{y}, t) - w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)| \, d\mathbf{x} \, dt \leq A + B,$$

where

$$\begin{aligned} A &= \int_0^{t_f} \int_{\mathbb{R}^2} |w_{\mathcal{T}_h, \Delta t_h}(\mathbf{x} + \mathbf{y}, t) - w_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t)| \, d\mathbf{x} \, dt, \\ B &= 2 \iint_{Q_{t_f}} |w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t) - w_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t)| \, d\mathbf{x} \, dt. \end{aligned}$$

One can conclude the proof using estimate (3.39) and the following estimate (resulting from a straightforward generalization of Lemma A.5 given in Appendix A or using Lemma 3.11)

$$\iint_{Q_{t_f}} |w_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t) - w_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t)| \, d\mathbf{x} \, dt \leq Ch. \quad (3.42)$$

This ends the proof of the lemma. \square

3.6 Convergence of the CVFE scheme

We can prove the main result of this section. Specifically, we have the following lemmas.

Lemma 3.11. *The sequences $(A(u_{\mathcal{M}_h, \Delta t_h}) - A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h}))_h$ and $(v_{\mathcal{M}_h, \Delta t_h} - v_{\mathcal{T}_h, \Delta t_h})_h$ converge strongly to zero in $L^2(Q_{t_f})$ as $h \rightarrow 0$.*

Proof. Using the definition of the basis functions of the finite dimensional space $\mathcal{H}_{\mathcal{T}_h}$, we have for all $M \in \mathcal{M}_h$ and all $K \in \mathcal{T}_h$ such that $M \cap K \neq \emptyset$

$$\begin{aligned} |A(u_{\mathcal{M}_h, \Delta t_h}) - A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h})|^2 &= |A(u_{\mathcal{M}_h, \Delta t_h})(\mathbf{x}_M, t_{n+1}) - A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h})(\mathbf{x}, t_{n+1})|^2 \\ &= |\nabla A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h})(\mathbf{x}, t_{n+1}) \cdot (\mathbf{x}_M - \mathbf{x})|^2, \quad \forall \mathbf{x} \in K \cap M \end{aligned}$$

where \mathbf{x}_M is the center of the dual control volume $M \in \mathcal{M}_h$.

Using estimate (3.35), one obtains

$$\begin{aligned} \|A(u_{\mathcal{M}_h, \Delta t_h}) - A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h})\|_{L^2(Q_{t_f})}^2 &= \sum_{n=0}^N \Delta t_h \sum_{K \in \mathcal{T}_h} \sum_{M, M \cap K \neq \emptyset} \int_{K \cap M} |\nabla A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h})(\mathbf{x}, t_{n+1}) \cdot (\mathbf{x}_M - \mathbf{x})|^2 \, d\mathbf{x} \\ &\leq h^2 \sum_{n=0}^N \Delta t_h \sum_{K \in \mathcal{T}_h} \sum_{M, M \cap K \neq \emptyset} |K \cap M| \left| \nabla A_{\mathcal{T}_h}(u_{\mathcal{T}_h}^{n+1}) \right|_K^2 \\ &\leq h^2 \sum_{n=0}^N \Delta t_h \|A_{\mathcal{T}_h}(u_{\mathcal{T}_h}^{n+1})\|_{\mathcal{H}_{\mathcal{T}_h}}^2 \leq Ch^2. \end{aligned}$$

As a consequence, we have $\|A(u_{\mathcal{M}_h, \Delta t_h}) - A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h})\|_{L^2(Q_{t_f})} \rightarrow 0$ as $h \rightarrow 0$. In the same manner, we prove that $\|v_{\mathcal{M}_h, \Delta t_h} - v_{\mathcal{T}_h, \Delta t_h}\|_{L^2(Q_{t_f})} \rightarrow 0$ as $h \rightarrow 0$. \square

Lemma 3.12. *(Convergence of the scheme). Under the shape regularity assumption (3.11) and the nonnegativity of the transmissibility coefficients assumption (3.26), there exists a sequence $(h_m)_{m \in \mathbb{N}}$, $h_m \rightarrow 0$ as $m \rightarrow \infty$, and functions u, v defined in Q_{t_f} such that $0 \leq u \leq 1$, both $A(u)$ and v belong to $L^2(0, T; H^1(\Omega))$, and*

$$A_{h_m}(u_{h_m}) \rightarrow A(u) \text{ and } v_{h_m} \rightarrow v \text{ a.e. in } Q_{t_f} \text{ and strongly in } L^p(Q_{t_f}) \text{ for all } p < +\infty$$

Proof. Let us set $\tilde{A}_{\mathcal{M}_h, \Delta t_h} := A(u_{\mathcal{M}_h, \Delta t_h})$ in Q_{t_f} and $\tilde{A}_{\mathcal{M}_h, \Delta t_h} := 0$ in $\mathbb{R}^3 \setminus Q_{t_f}$. Thanks to Proposition 3.6 and Lemma 3.5, one has $(\tilde{A}_{\mathcal{M}_h, \Delta t_h}) \subset L^\infty(\mathbb{R}^3) \cap L^2(\mathbb{R}^3)$. In order to verify the assumptions of Kolmogorov's compactness criterion, see [31, theorem 3.9, p. 93], we note that the following inequality is verified for any $\eta \in \mathbb{R}^2$ and $\tau \in \mathbb{R}$,

$$\begin{aligned} \left\| \tilde{A}_{\mathcal{M}_h, \Delta t_h}(\cdot + \eta, \cdot + \tau) - \tilde{A}_{\mathcal{M}_h, \Delta t_h}(\cdot, \cdot) \right\|_{L^1(\mathbb{R}^3)} &\leq \left\| \tilde{A}_{\mathcal{M}_h, \Delta t_h}(\cdot + \eta, \cdot) - \tilde{A}_{\mathcal{M}_h, \Delta t_h}(\cdot, \cdot) \right\|_{L^1(\mathbb{R}^3)} \\ &\quad + \left\| \tilde{A}_{\mathcal{M}_h, \Delta t_h}(\cdot, \cdot + \tau) - \tilde{A}_{\mathcal{M}_h, \Delta t_h}(\cdot, \cdot) \right\|_{L^1(\mathbb{R}^3)}. \end{aligned}$$

Using estimates (3.36) and (3.38), one has $\left\| \tilde{A}_{\mathcal{M}_h, \Delta t_h}(\cdot + \eta, \cdot + \tau) - \tilde{A}_{\mathcal{M}_h, \Delta t_h}(\cdot, \cdot) \right\|_{L^1(\mathbb{R}^3)} \rightarrow 0$, as $\eta \rightarrow 0$ and $\tau \rightarrow 0$. This yields the compactness of the sequence $(\tilde{A}_{\mathcal{M}_h, \Delta t_h})$ in $L^1(\Omega)$.

Thus, there exists a subsequence, still denoted by $(\tilde{A}_{\mathcal{M}_h, \Delta t_h})$, and there exists $A^* \in L^1(Q_{t_f})$ such that

$$A(u_{\mathcal{M}_h, \Delta t_h}) \rightarrow A^* \text{ strongly in } L^1(Q_{t_f}).$$

Furthermore, as A is continuous and strictly monotone, there exists a unique u such that $A(u) = A^*$.

Since A^{-1} is well defined and continuous, then applying the L^∞ bound on $u_{\mathcal{M}_h, \Delta t_h}$ and the dominated convergence theorem to $u_{\mathcal{M}_h, \Delta t_h} = A^{-1}(A(u_{\mathcal{M}_h, \Delta t_h}))$, we get

$$u_{\mathcal{M}_h, \Delta t_h} \rightarrow u \text{ a.e. in } Q_{t_f} \text{ and strongly in } L^p(Q_{t_f}) \text{ for } p < +\infty.$$

According to Lemma 3.11, the sequences $(A(u_{\mathcal{M}_h, \Delta t_h}))_{h, \Delta t_h}$ and $(A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h}))_{h, \Delta t_h}$ have the same limit, as a consequence

$$A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h}) \rightarrow A(u) \text{ strongly in } L^2(Q_{t_f}) \text{ and a.e. in } Q_{t_f}.$$

Similarly, translate estimates (3.36)–(3.38), the L^∞ bound on $v_{h, \Delta t_h}$ in Lemma 3.5, and Lemma 3.11 ensure that, up to extraction of a subsequence,

$$v_{\mathcal{T}_h, \Delta t_h} \rightarrow v \text{ a.e. in } Q_{t_f} \text{ and strongly in } L^p(Q_{t_f}) \text{ for } p < +\infty.$$

It remains to show that $A(u)$ and $v \in L^2(0, T; H^1(\Omega))$. Indeed, using estimate (3.33), one gets that the sequence $\nabla A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h}) \in (L^2(Q_{t_f}))^2$. It follows that $(A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h}))$ is bounded in $L^2(0, T; H^1(\Omega))$ since $A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h})$ is uniformly bounded in $L^2(Q_{t_f})$ due to the boundedness of the function A . Therefore, the sequence $(A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h}))$ converges weakly, up to an unlabeled subsequence, to a function A^* in $L^2(0, T; H^1(\Omega))$. The sequence $(A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h}))$ converges strongly in $L^2(Q_{t_f})$ to $A(u)$, one can deduce that $A(u) = A^* \in L^2(0, T; H^1(\Omega))$. In the same manner, we obtain that $v \in L^2(0, T; H^1(\Omega))$. This ends the proof of the lemma. \square

It remains to be shown that the limit functions u and v constitute a weak solution of the continuous system. For this, let $\psi \in \mathcal{D}([0, T] \times \bar{\Omega})$ be a test function and denote by $\psi_M^n :=$

$\psi(\mathbf{x}_M, t_n)$ for all $M \in \mathcal{M}_h$ and $n \in \{0, \dots, N+1\}$.

Multiply equation (3.24) by $\Delta t_h \psi_M^{n+1}$ and sum up over $M \in \mathcal{M}_h$ and $n \in \{0, \dots, N\}$. One has

$$\begin{aligned} & \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} |M| (u_M^{n+1} - u_M^n) \psi_M^{n+1} \\ & + \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (A(u_M^{n+1}) - A(u_{M'}^{n+1})) \psi_M^{n+1} \\ & + \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) |\sigma_{M,M'}^K| \psi_M^{n+1} \\ & = \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} |M| f(u_M^{n+1}) \psi_M^{n+1}, \end{aligned}$$

that is,

$$S_1^h + S_2^h + S_3^h = S_4^h$$

where

$$\begin{aligned} S_1^h &:= \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} |M| (u_M^{n+1} - u_M^n) \psi_M^{n+1}, \quad S_4^h := \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} |M| f(u_M^{n+1}) \psi_M^{n+1}, \\ S_2^h &:= \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^K) \psi_M^{n+1} |\sigma_{M,M'}^K|, \\ S_3^h &:= \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (A(u_M^{n+1}) - A(u_{M'}^{n+1})) \psi_M^{n+1}. \end{aligned}$$

Perform a summation by parts in time and keep in mind that $\psi_M^{N+1} = 0$ for all $M \in \mathcal{M}_h$, we obtain

$$\begin{aligned} S_1^h &= \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} |M| u_M^{n+1} \psi_M^{n+1} - \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} |M| u_M^n \psi_M^{n+1} \\ &= - \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} \Delta t_h |M| u_M^n \frac{\psi_M^{n+1} - \psi_M^n}{\Delta t_h} - \sum_{M \in \mathcal{M}} |M| u_M^0 \psi_M^0 \\ &= - \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} \int_{t_n}^{t_{n+1}} \int_M u_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t) \partial_t \psi(\mathbf{x}_M, t) d\mathbf{x} dt \\ &\quad - \sum_{M \in \mathcal{M}} \int_M u_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, 0) \psi(\mathbf{x}_M, 0) d\mathbf{x}. \end{aligned}$$

Taking into account the assumptions on the data and using the Lebesgue theorem, it follows that

$$S_1^h \xrightarrow{h, \Delta t_h \rightarrow 0} - \int_0^{t_f} \int_{\Omega} u \partial_t \psi d\mathbf{x} dt - \int_{\Omega} u_0 \psi(\cdot, 0) d\mathbf{x}.$$

On the other hand, for the convergence of the third term S_3^h , we note that

$$\begin{aligned}
S_3^h &= \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} (\Lambda_{M,M'}^K A(u_M^{n+1}) \psi_M^{n+1} - \Lambda_{M,M'}^K A(u_{M'}^{n+1}) \psi_{M'}^{n+1}) \\
&= - \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} A(u_M^{n+1}) \sum_{\sigma_{M,M'}^K \subset \partial M} \Lambda_{M,M'}^K (\psi_{M'}^{n+1} - \psi_M^{n+1}) \\
&= - \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} \Delta t_h A(u_M^{n+1}) \sum_{K \cap M \neq \emptyset} \sum_{\sigma_{M,M'}^K \subset \partial M \cap K} \nabla \psi^{n+1}|_K \cdot {}^t \Lambda_K \eta_{M,\sigma} |\sigma| \\
&= - \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} \Delta t_h A(u_M^{n+1}) \sum_{\sigma_{M,M'}^K \subset \partial M} \nabla \psi^{n+1}|_K \cdot {}^t \Lambda_K \eta_{M,\sigma} |\sigma| \\
&= - \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} \Delta t_h A(u_M^{n+1}) \int_{\partial M} \Lambda \nabla \psi^{n+1} \cdot \eta \, d\sigma(\mathbf{x}) \\
&= - \sum_{n=0}^N \sum_{M \in \mathcal{M}_h} \int_{t_n}^{t_{n+1}} \int_M A(u_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)) \operatorname{div}(\Lambda \nabla \psi) \, d\mathbf{x} \, dt \\
&\xrightarrow{h, \Delta t_h \rightarrow 0} - \iint_{Q_{t_f}} A(u) \operatorname{div}(\Lambda \nabla \psi) \, d\mathbf{x} \, dt = \iint_{Q_{t_f}} \nabla A(u) \cdot \Lambda \nabla \psi \, d\mathbf{x} \, dt.
\end{aligned}$$

It remains to show that

$$\lim_{h, \Delta t_h \rightarrow \infty} S_2^h = - \int_0^{t_f} \int_{\Omega} \Lambda(\mathbf{x}) \chi(u) \nabla v \cdot \nabla \psi \, d\mathbf{x} \, dt. \quad (3.43)$$

For the convergence of S_2^h , we note that gathering by edges (thanks to the consistency of the fluxes, see (3.20)(b)), we find

$$S_2^h = -\frac{1}{2} \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M,M'}^K \subset \partial M} G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M,M'}^{n+1}) (\psi_{M'}^{n+1} - \psi_M^{n+1}) |\sigma_{M,M'}^K|.$$

For each triplet of neighbors M , M' , and M'' pick for $u_{K,\min}^{n+1}$ the quantity defined by

$$u_{K,\min}^{n+1} = \min_{M \in \mathcal{M}, M \cap K \neq \emptyset} \{u_M^{n+1}\}$$

Set

$$S_2^{h,*} := \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} \sum_{K, K \cap M \neq \emptyset} \sum_{M'', M'' \cap K \neq \emptyset} \chi(u_{K,\min}^{n+1}) \psi_M^{n+1} dV_{M,M'}^{n+1} |\sigma_{M,M'}^K|.$$

We have

$$\begin{aligned}
S_2^{h,*} &= - \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} \sum_{K, K \cap M \neq \emptyset} \chi(u_{K,\min}^{n+1}) |K| \Lambda_K \psi_M^{n+1} \nabla \varphi_M \cdot \left(\sum_{M'', M'' \cap K \neq \emptyset} v_{M''}^{n+1} \nabla \varphi_{M''} \right) \\
&= - \sum_{n=0}^N \Delta t_h \sum_{K \in \mathcal{T}_h} |K| \chi(u_{K,\min}^{n+1}) \Lambda_K \left(\sum_{M, M \cap K \neq \emptyset} \psi_M^{n+1} \nabla \varphi_M \right) \cdot \left(\sum_{M'', M'' \cap K \neq \emptyset} v_{M''}^{n+1} \nabla \varphi_{M''} \right).
\end{aligned}$$

Introduce $\bar{u}_h, \underline{u}_h$ defined by

$$\bar{u}_h|_{(t_n, t_{n+1}] \times K} := \max_{M, M \cap K \neq \emptyset} \{u_M^{n+1}\}, \quad \underline{u}_h|_{(t_n, t_{n+1}] \times K} := \min_{M, M \cap K \neq \emptyset} \{u_M^{n+1}\}.$$

Consequently, we obtain

$$\begin{aligned} S_2^{h,*} &= - \sum_{n=0}^N \Delta t_h \sum_{K \in \mathcal{T}_h} \int_K \chi(\underline{u}_h) \Lambda_K \nabla v_{\mathcal{T}_h, \Delta t_h}^{n+1} \cdot (\nabla \psi)_{\mathcal{T}_h}^{n+1} d\mathbf{x} \\ &= - \int_{Q_{t_f}} \chi(\underline{u}_h) \Lambda \nabla v_{h, \Delta t_h} \cdot \nabla \psi_{\mathcal{T}_h, \Delta t_h} d\mathbf{x} dt. \end{aligned}$$

Next, we show that

$$\lim_{h \rightarrow 0} |S_2^h - S_2^{h,*}| = 0. \quad (3.44)$$

To do this, we begin by showing that $|\underline{u}_h - \bar{u}_h| \rightarrow 0$ a.e. in Q_{t_f} .

By a slight adaptation of Lemma A.1 and thanks to estimate (3.33), one can conclude the existence of a constant $C > 0$ (independent of h) such that

$$\begin{aligned} \int_0^{t_f} \int_{\Omega} |A(\bar{u}_h) - A(\underline{u}_h)|^2 d\mathbf{x} dt &\leq C h^2 \int_0^{t_f} \int_{\Omega} |\nabla A_{\mathcal{T}_h}(u_{\mathcal{T}_h, \Delta t_h})(\mathbf{x}, t)|^2 d\mathbf{x} dt \\ &\leq C \Lambda_- h^2 \sum_{n=0}^N \Delta t_h \sum_{M \in \mathcal{M}_h} \sum_{\sigma_{M, M'}^K \subset \partial M} \Lambda_{M, M'}^K |A(u_{M'}^{n+1}) - A(u_M^{n+1})|^2 \leq C h^2. \end{aligned}$$

Since A^{-1} is continuous, up to extraction of another subsequence, we deduce

$$|\bar{u}_h - \underline{u}_h| \rightarrow 0 \text{ a.e. in } Q_{t_f}. \quad (3.45)$$

In addition, $\underline{u}_h \leq u_{\mathcal{M}_h, \Delta t_h} \leq \bar{u}_h$; moreover, by Lemma 3.12, $u_{\mathcal{M}_h, \Delta t_h} \rightarrow u$ a.e. in Q_{t_f} . Thus we see that $\chi(\underline{u}_h) \rightarrow \chi(u)$ a.e. in Q_{t_f} and in $L^p(Q_{t_f})$, for $p < +\infty$. Using Lemma 3.12 again, the strong convergence of $\nabla \psi_{\mathcal{T}_h, \Delta t_h}$ towards $\nabla \psi$, and the weak convergence in $L^2(Q_{t_f})$ of $\nabla v_{\mathcal{T}_h, \Delta t_h}$ towards ∇v , we conclude that

$$\lim_{h \rightarrow 0} S_2^{h,*} = - \int_0^{t_f} \int_{\Omega} \chi(u) \Lambda \nabla v \cdot \nabla \psi d\mathbf{x} dt.$$

To prove (3.44), we remark that

$$\begin{aligned} &\left| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M, M'}^{n+1}) - \chi(u_{K, \min}^{n+1}) dV_{M, M'}^{n+1} \right| \\ &= \left| G(u_M^{n+1}, u_{M'}^{n+1}; dV_{M, M'}^{n+1}) - G(u_{K, \min}^{n+1}, u_{K, \min}^{n+1}; dV_{M, M'}^{n+1}) \right| \\ &\leq \left| dV_{M, M'}^{n+1} \right| \omega(2|u_M^{n+1} - u_{K, \min}^{n+1}|). \end{aligned}$$

Consequently, we obtain

$$|S_2^h - S_2^{h,*}| \leq \int_0^{t_f} \int_{\Omega} \omega(2|\bar{u}_h - \underline{u}_h|) |\nabla v_{\mathcal{T}_h, \Delta t_h} \cdot \nabla \psi_{\mathcal{T}_h, \Delta t_h}| d\mathbf{x} dt.$$

Applying the Cauchy-Schwarz inequality, and the convergence (3.44), we establish (3.43).

Finally, we note that it is easy to see that $S_4^h \xrightarrow{h, \Delta t_h \rightarrow 0} \int_{Q_{t_f}} f(u) \psi d\mathbf{x} dt$.

3.7 Numerical simulation in two-dimensional space

In this section, we exhibit various two-dimensional numerical results provided by scheme (3.23)–(3.25) for the capture of spatial patterns for model (3.1) discussed in Section 3.2. Newton's algorithm is used to approach the solution U^{n+1} of the nonlinear system defined by equation (3.24), this algorithm is coupled with a bigradient method to solve the linear systems arising from the Newton algorithm as well as the linear system given by equation (3.25). Unless stated otherwise, throughout this section, we consider that the cell density is initially set as a spatially small random perturbation around the homogeneous steady state, and we assume zero-flux boundary conditions. The simulations are performed on an unstructured triangular mesh of the space domain $\Omega = (0, 10) \times (0, 10)$. We suppose that the species cells follow the logistic growth $f(u) = \mu u(1 - u/u_c)$, where μ is the *intrinsic growth rate*, and u_c is the *carrying capacity* of the population. Production term $g(u, v)$, squeezing probability $q(u)$, cell diffusivity $a(u)$, and chemotactic sensitivity $\chi(u)$ are given by (3.4), (3.7), and (3.8) respectively.

The pattern formation for model (3.1) with the associated functions mentioned above has been established in [78] using Turing's principle and the linear stability analysis, where the diffusion tensor is considered to be proportional to the identity matrix and the numerical simulations are carried on a one-dimensional space, while in [49], the same analysis is provided whereas the numerical simulations for the capture of spatial patterns are presented on a two-dimensional domain, and using the standard finite volume scheme.

The nontrivial uniform steady state of system (3.1) is given by $(u_s, v_s) = (u_c, \alpha u_c / \beta)$, and through the pattern formation analysis provided in [78, 49], the instability region of this steady state is determined by the following condition :

$$\mu + \beta a(u_c) - \alpha \chi(u_c) < -2\sqrt{\mu \beta a(u_c)}. \quad (3.46)$$

In order to verify the effectiveness of the proposed scheme, we consider three tests for the capture of spatio-temporal patterns for model (3.1) with different diffusion tensors. For each test, we choose a set of parameters, such that the instability condition (3.46) is satisfied. We fix $d_1 = 0.25$, $u_c = 0.25$, $\bar{u} = 1.0$, $\mu = 0.5$, $\alpha = 10.0$, $\beta = 10.0$, $\zeta = 20$, and $\gamma = 3$. On the other hand, and in the definition of the numerical flux function G defined by (3.20), we take

$$\chi_{\uparrow}(z) = \chi(\min\{z, \tilde{u}\}) \text{ and } \chi_{\downarrow}(z) = \chi(\max\{z, \tilde{u}\}) - \chi(\tilde{u}),$$

where $\tilde{u} = \frac{\bar{u}}{\sqrt[\gamma]{\gamma+1}}$. Finally, we consider a small time step $\Delta t = 0.005$ and a nonuniform primary mesh with small refinement consisting of 14 336 triangles. Thus, the associated Donald dual mesh consists of 7297 dual control volumes.

Test 1 (Isotropic case) Here, we establish the generation of spatial patterns for the volume-filling chemotaxis model (3.1) with homogeneous diffusion tensors. In this test, we take $\Lambda(\mathbf{x}) = D(\mathbf{x}) = I_2(\mathbf{x})$ and initially the chemical concentration is set to be a constant equal to v_s .

Figures 3.2–3.3 show for different moments, the pattern formation for model (3.1) with identity diffusion tensors. We see that the random distribution of the cell density leads to a merging process in all directions of the space at $t = 0.55s$, which continues for $t = 2.35s$ then it stops when the time $t \geq 20.75s$, and new stationary spot patterns appear as shown at $t = 40s$.

Time evolution of the cell density. Here we consider the time evolution of the cell density at fixed points in the right snapshot of Figure 3.3. Indeed, we want to show that the cell density stabilizes at a certain moment; hence, we prove that the volume-filling chemotaxis model (3.1) generates stationary spatial patterns. Figure 3.4 shows the evolution of the cell density with respect

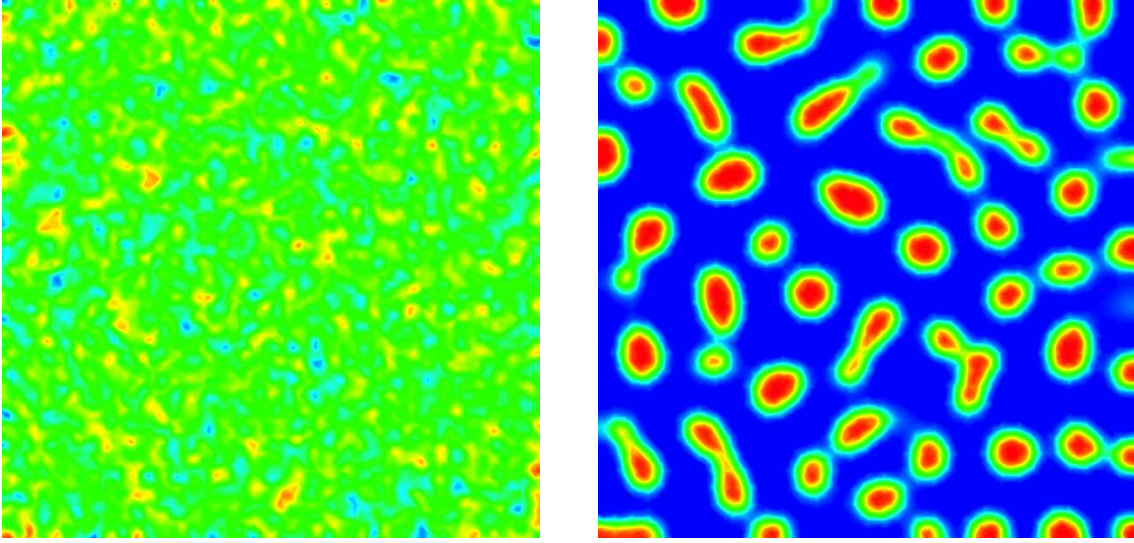


FIGURE 3.2 – Pattern formation for the full chemotaxis model (3.1) on a 2-D domain $\Omega = (0, 10) \times (0, 10)$ at time $t = 0$ s with $0 \leq u \leq 1$ (*left*) and at time $t = 0.55$ s with $3.10^{-3} \leq u \leq 0.99$ (*right*).

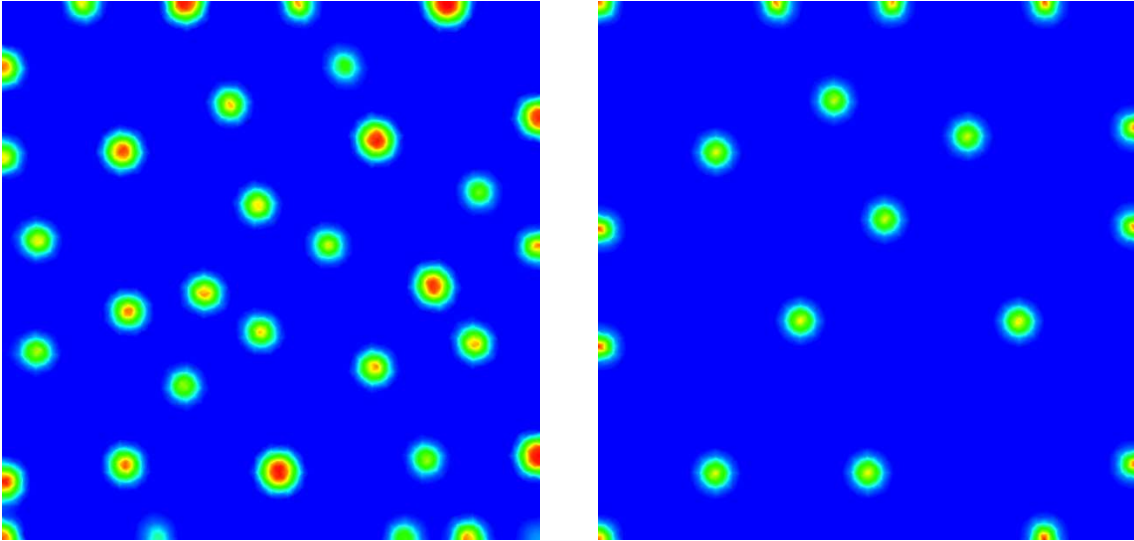


FIGURE 3.3 – Pattern formation for the full chemotaxis model (3.1) at time $t = 2.35$ s with $1.9 \times 10^{-4} \leq u \leq 0.98$ (*left*) and at time $t = 40$ s with $3.9 \times 10^{-3} \leq u \leq 0.832$ (*right*). The red dots (or rods) represent the cell aggregation where cell density is higher than that of the blue area.

to the time at point $P_1 (5.25; 6)$ in the red line, at point $P_2 (5; 1.25)$ in the green line, and at point $P_3 (4.35; 8.1)$ in the blue line. We observe that the cell density at these points increases and then decreases with response to the gradient of the chemoattractant which plays an essential role to stop the aggregation of the cells. Next, the cell density stabilizes for all points when t is greater or equal to 13 s. We note that the same results are obtained for the other spot patterns; however, for the sake of brevity, they are not provided here.

Test 2 (Anisotropic case) In this test, we take a homogeneous anisotropic diffusion tensors of the form $\begin{bmatrix} 1 & 0 \\ 0 & \xi \end{bmatrix}$. We investigate the pattern formation for model (3.1) and consider that the

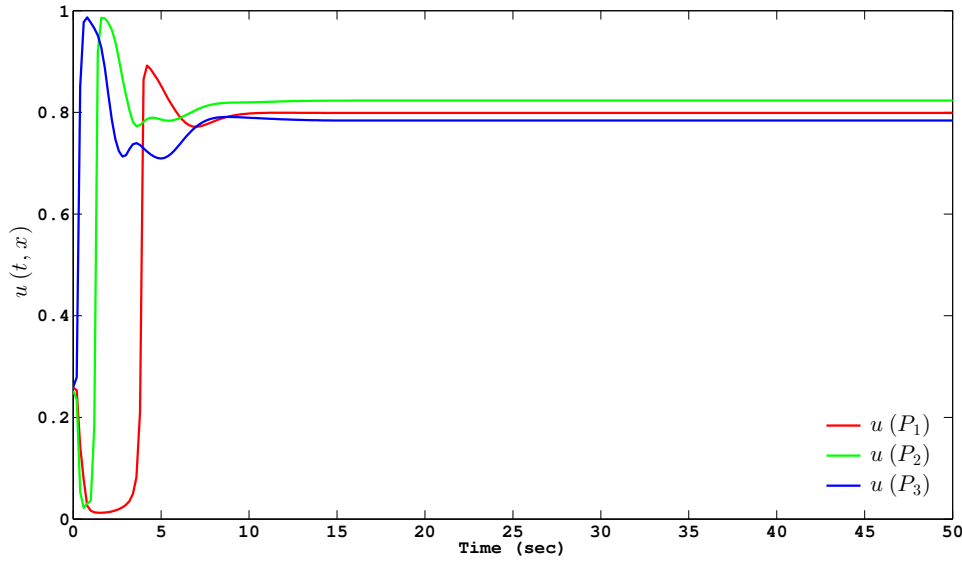


FIGURE 3.4 – Evolution of the cell density with respect to the time at three different points : P_1 (5.25; 6) in the red line, P_2 (5; 1.25) in the green line, P_3 (4.35; 8.1) in the blue line.

particles diffuse more rapidly in the x -axis direction than the y -axis direction. In other words, we are concerned by a matrix with high diffusivity in the horizontal direction and low diffusivity in the orthogonal direction ; for that we take $\xi < 1$, say for example $\xi = 0.5$. We pick up snapshots for the pattern formation at the same moments as those of the isotropic case.

Figures 3.5–3.6 show the evolution of spatial patterns for model (3.1) for different time moments. We observe that we have the same results as before (same patterning, and emerging process), except that more spot patterns are obtained and they are stretched in the horizontal direction as shown in the last snapshot in figure 3.6.

The time evolution for the last plot in figure 3.6 for different spots is given in figure 3.7. It shows that the stationary spot patterns are obtained when $t \geq 23.5s$.

Test 3 (Heterogeneous anisotropic case) In this test, we decompose the domain Ω into two regions Ω_1 and Ω_2 , where $\Omega_1 = (0, 10) \times (0, 5]$ and $\Omega_2 = (0, 10) \times (5, 10)$. Moreover, we assume that the diffusion tensors are anisotropic and heterogeneous, and are given by :

$$\Lambda(\mathbf{x}) = D(\mathbf{x}) = \begin{pmatrix} 1 & 0 \\ 0 & \lambda(\mathbf{x}) \end{pmatrix}, \text{ with } \begin{cases} \lambda(\mathbf{x}) = 0.5, & \text{if } \mathbf{x} \in \Omega_1, \\ \lambda(\mathbf{x}) = 1.5, & \text{if } \mathbf{x} \in \Omega_2. \end{cases}$$

Figures 3.8–3.9 show the evolution of spatial patterns for model (3.1) for different time moments.

In region Ω_2 , where diffusion is more interesting in the y -axis direction, and at $t = 0.5s$, we see that the aggregations form quickly, and they are much larger than those formed in region Ω_1 , which means that high diffusion in region Ω_2 (compared to that in region Ω_1) accelerate the merging process. On the other hand, at $t = 49s$ the aggregations in region Ω_2 disappear to form 5 spatial patterns ; whereas, a plentitude of spatial patterns is seen in region Ω_1 compared to region Ω_2 . Therefore, there exists a high dependence between the rate of diffusion and the generation of spatial patterns.

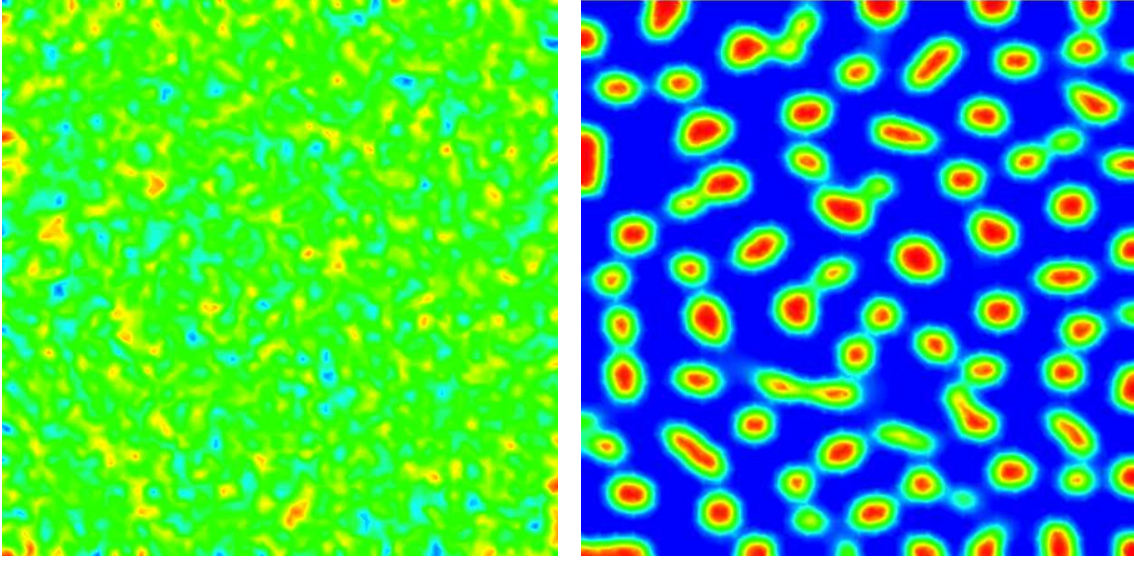


FIGURE 3.5 – Pattern formation for the full chemotaxis model (3.1) on a 2-D domain $\Omega = (0, 10) \times (0, 10)$ at time $t = 0$ s with $0 \leq u \leq 1$ (*left*) and at time $t = 0.55$ s with $4.10^{-3} \leq u \leq 0.98$ (*right*).

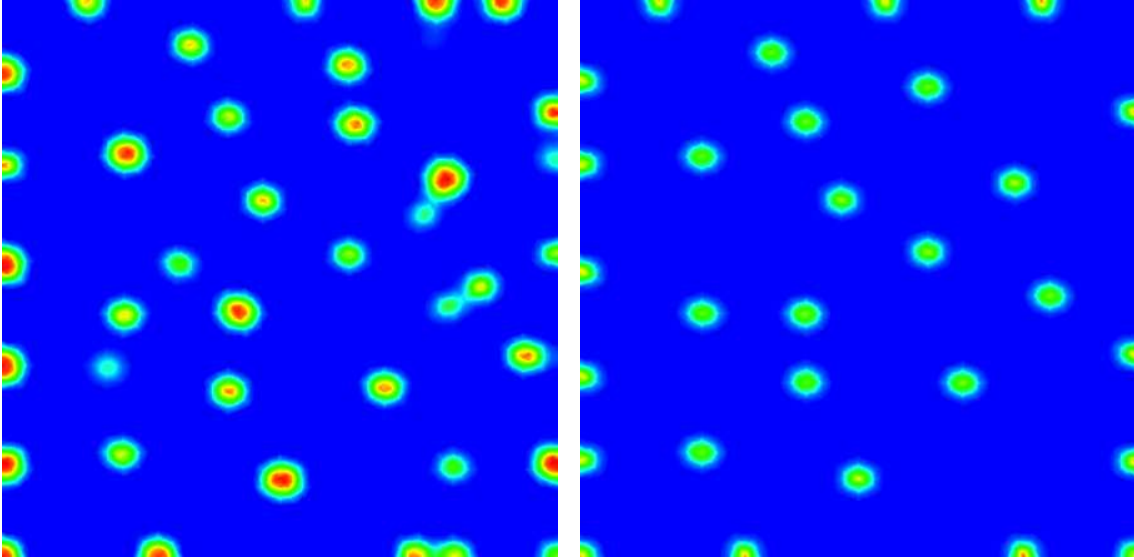


FIGURE 3.6 – Pattern formation for the full chemotaxis model (3.1) at time $t = 2.35$ s with $3.2 \times 10^{-4} \leq u \leq 0.98$ (*left*) and at time $t = 40$ s with $4 \times 10^{-3} \leq u \leq 0.827$ (*right*).

The time evolution for the last plot of figure 3.9 for different spots is given in figure 3.10. It shows that stationary spot patterns are obtained when $t \geq 125$ s.

Comparing figures 3.5-3.9, one can deduce that the patterning depends on diffusion coefficients, it is also known to depend on the size of the domain. These results prove the robustness of the *control volume finite element scheme* to capture spatial patterns for a volume-filling chemotaxis model with anisotropic diffusion tensors.

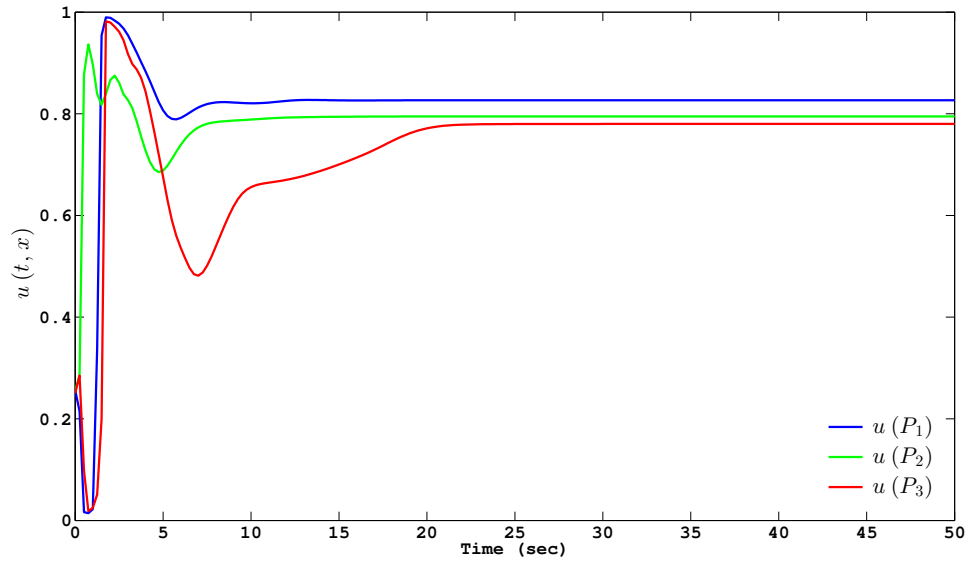


FIGURE 3.7 – The evolution of the cell density with respect to the time at three different points : P_1 (5; 1.4) in the blue line, P_2 (6.25; 5.45) in the green line, P_3 (9.9; 8.05) in the red line.

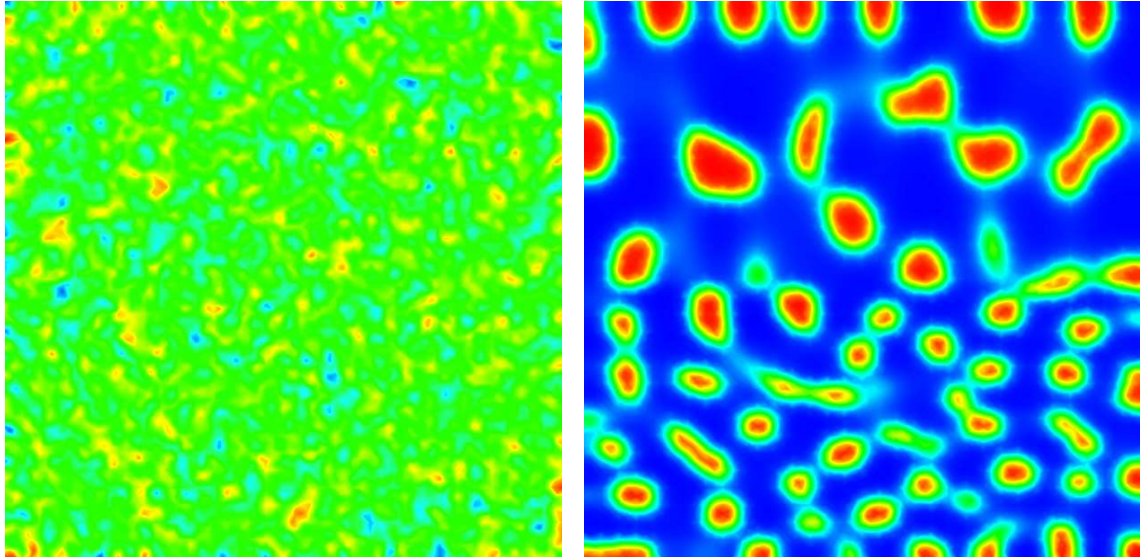


FIGURE 3.8 – Pattern formation for the full chemotaxis model (3.1) on a 2-D domain $\Omega = (0, 10) \times (0, 10)$ at time $t = 0$ s with $0 \leq u \leq 1$ (left) and at time $t = 0.5$ s with $6.93 \times 10^{-3} \leq u \leq 0.99$ (right).

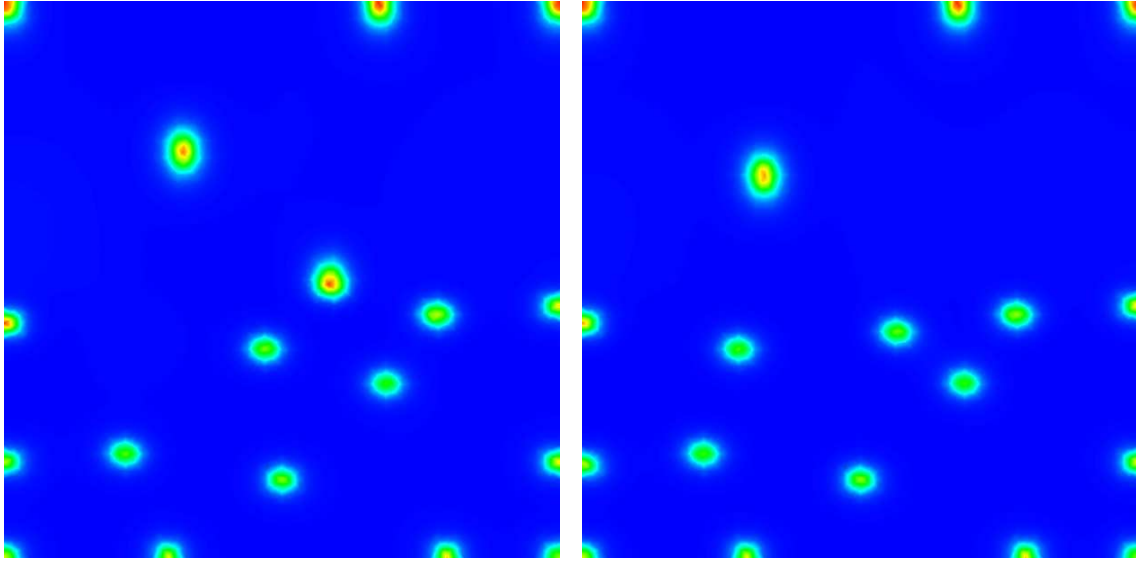


FIGURE 3.9 – Pattern formation for the full chemotaxis model (3.1) at time $t = 49\text{s}$ with $3.67 \times 10^{-3} \leq u \leq 0.88$ (*left*) and at time $t = 150\text{s}$ with $4.3 \times 10^{-3} \leq u \leq 0.8447$ (*right*).

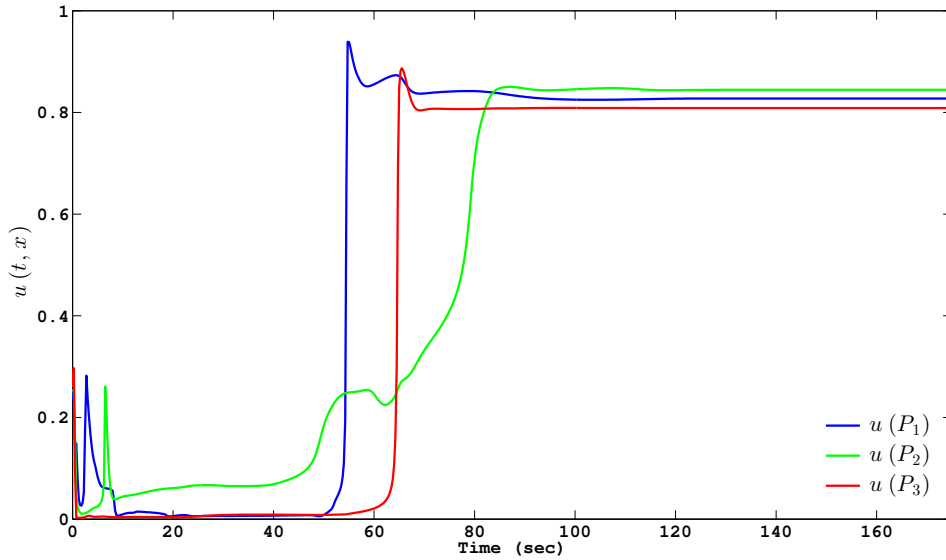


FIGURE 3.10 – The evolution of the cell density with time at three different points : P_1 (5.70; 4.05) in the blue line, P_2 (3.28; 6.88) in the green line, P_3 (2.80; 3.75) in the red line.

A nonlinear CVFE scheme for the modified degenerate anisotropic chemotaxis model

Sommaire

4.1	Introduction	86
4.2	The modified chemotaxis model	86
4.3	Space-time discretization and notations	88
4.3.1	Space discretizations of Ω .	88
4.3.2	Discrete finite elements space $\mathcal{H}_{\mathcal{T}}$, control volumes space $\mathcal{X}_{\mathcal{M}}$.	88
4.3.3	Time discretization of $(0, t_f)$.	89
4.3.4	Space-time discretization of Q_{t_f} .	89
4.3.5	Main property	90
4.4	The nonlinear CVFE scheme	90
4.4.1	Discretization of the first equation of system (4.1)	91
4.4.2	Discretization of the second equation of system (4.1)	93
4.4.3	Main result	94
4.5	Discrete properties, a priori estimates and existence	94
4.5.1	Discrete maximum principle	97
4.5.2	Entropy estimates on $v_{\mathcal{M}, \Delta t}$	99
4.5.3	Energy estimates on $u_{\mathcal{M}, \Delta t}$	101
4.5.4	Enhanced estimate on $v_{\mathcal{M}, \Delta t}$	103
4.5.5	Existence of a discrete solution	104
4.6	Compactness estimates on the family of discrete solutions.	106
4.6.1	Time translate estimate.	106
4.6.2	Space translate estimate.	108
4.7	Convergence	109
4.7.1	Identification as a weak solution	110
4.8	Numerical results	114

4.1 Introduction

In this chapter, we are interested in a degenerate nonlinear parabolic reaction–convection–diffusion system modeling the chemotaxis process over general triangular mesh, with anisotropic and heterogeneous diffusion tensors.

In chapter 3, the convergence analysis of the mixed conforming piecewise linear finite elements on triangles for the diffusion term and finite volume on dual elements for the other terms has carried out for the case of anisotropic and heterogeneous diffusion problems under an essential assumption that all the transmissibility coefficients are nonnegative. However, there is no sufficient conditions for nonnegativity of transmissibility coefficients and therefore the schemes do not permit to tackle general anisotropic diffusion problems.

Recently, Cancès and Guichard proposed and analyzed in [16] a nonlinear *Control Volume Finite Element* (CVFE) scheme for solving degenerate anisotropic parabolic diffusion equations modeling flows in porous media. The convergence analysis is carried out without any restriction on the transmissibility coefficients, and the effectiveness of the scheme is tested using anisotropic diffusion tensors over an unstructured mesh.

The purpose of this chapter is to extend the idea of [16] to a fully nonlinear degenerate parabolic reaction–convection–diffusion system modeling the chemotaxis process over general mesh, with anisotropic and heterogeneous diffusion tensors. Here, the applied discretization is very similar to that one employed in [50]; the diffusion terms are discretized by means of a conforming piecewise linear finite element method on a primal triangular mesh and using the Godunov scheme to approximate the diffusion fluxes provided by the conforming finite element reconstruction. The others terms are discretized by means of a nonclassical upwind finite volume method on a dual mesh (Donald mesh or Median dual mesh).

The rest of this chapter is organized as follows. In Section 4.2, we introduce the modified Keller-Segel model as well as the assumptions made about the system. In Section 4.3 and in order to discretize the domain, we define a primal triangular mesh and its corresponding dual barycentric mesh, we then define two different reconstructions of the discrete solutions based on the spatial discretizations; the first stands to the usual \mathbb{P}_1 finite element reconstruction while the second stands to the piecewise constant reconstruction. In Section 4.4, we define the discretization of the diffusion and convection terms and introduce the nonlinear CVFE scheme. With the help of some technical lemmas in Section 4.5, we prove the existence of a discrete solution to the CVFE scheme based on the establishment of *a priori* estimates on the discrete solution as well as the discrete maximum principle. In Section 4.6, we give estimates on differences of time and space translates for the approximate solutions. In Section 4.7, using the Kolmogorov relative compactness criterion, we prove the convergence of a subsequent of the sequence of discrete solutions to a weak solution of the continuous problem. Finally, some numerical simulations are carried out, in Section 4.8, to show the effectiveness of the scheme to tackle degenerate anisotropic and heterogeneous diffusion problems over general unstructured mesh.

4.2 The modified chemotaxis model

Let Ω be an open bounded connected polygonal domain of \mathbb{R}^2 , and $t_f > 0$ be a fixed time. The modified Keller-Segel system modeling the chemotaxis process is given by the following set of equations

$$\begin{cases} \partial_t u - \operatorname{div} (\Lambda(\mathbf{x}) a(u) \nabla u - \Lambda(\mathbf{x}) \chi(u) \nabla v) = f(u) & \text{in } Q_{t_f} = \Omega \times (0, t_f), \\ \partial_t v - \operatorname{div} (D(\mathbf{x}) \nabla v) = g(u, v) & \text{in } Q_{t_f} = \Omega \times (0, t_f). \end{cases} \quad (4.1)$$

The system is complemented with homogeneous zeros-flux boundary conditions on $\Sigma_{t_f} := \partial\Omega \times (0, t_f)$ given by

$$(\Lambda(\mathbf{x}) a(u) \nabla u - \Lambda(\mathbf{x}) \chi(u) \nabla v) \cdot \mathbf{n} = 0, \quad D(\mathbf{x}) \nabla v \cdot \mathbf{n} = 0, \quad (4.2)$$

and the initial conditions on Ω :

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad v(\mathbf{x}, 0) = v_0(\mathbf{x}). \quad (4.3)$$

In the above model, the density of the cell-population and the chemoattractant concentration are represented by $u = u(\mathbf{x}, t)$ and $v = v(\mathbf{x}, t)$ respectively. Next, $a(u)$ is a density-dependent diffusion coefficient, and $\Lambda(\mathbf{x})$ is the diffusion tensor in a heterogeneous medium. Furthermore, the function $\chi(u)$ is the chemoattractant sensitivity, and $D(\mathbf{x})$ is the diffusion tensor for v . The function $f(u)$ describes cell proliferation and cell death. The function $g(u, v)$ describes the rates of production and degradation of the chemoattractant ; here, we assume it is the linear function given by

$$g(u, v) = \alpha u - \beta v, \quad \alpha, \beta \geq 0. \quad (4.4)$$

The unit normal vector to $\partial\Omega$ outward to Ω is denoted by \mathbf{n} .

Let us state the main assumptions made about model (4.1)–(4.3) :

- (A1) The cell-density diffusion $a : [0, 1] \rightarrow \mathbb{R}^+$ is a continuous function such that, $a(0) = a(1) = 0$, and $a(u) > 0$ for $0 < u < 1$.
- (A2) The chemosensitivity $\chi : [0, 1] \rightarrow \mathbb{R}^+$ is a continuous function such that, $\chi(0) = \chi(1) = 0$. Furthermore, we assume that there exists a function $\mu \in C([0, 1]; \mathbb{R}^+)$, such that $\mu(u) = \frac{\chi(u)}{a(u)}$ and $\mu(0) = \mu(1) = 0$.
- (A3) The diffusion tensors Λ and D are two bounded, uniformly positive symmetric tensors on Ω , that is : $\forall \mathbf{w} \neq 0$, there exists two nonnegative constants T_- and T_+ such that $0 < T_- |\mathbf{w}|^2 \leq \langle T(\mathbf{x}) \mathbf{w}, \mathbf{w} \rangle \leq T_+ |\mathbf{w}|^2 < \infty$, $T = \Lambda$ or D .
- (A4) The cell density proliferation f is a continuous function such that $f(0) \geq 0$ and $f(1) \leq 0$.
- (A5) The initial function u_0 and v_0 are two functions in $L^2(\Omega)$ such that, $0 \leq u_0 \leq 1$ and $v_0 \geq 0$.

In the sequel, we use the Lipschitz continuous nondecreasing function $\xi : [0, 1] \rightarrow \mathbb{R}$ defined by

$$\xi(u) := \int_0^u \sqrt{a(s)} \, ds, \quad \forall u \in \mathbb{R}. \quad (4.5)$$

We recall the definition of a weak solution of system (4.1)–(4.3).

Definition 4.1 (weak solution). Under the assumptions (A1)–(A5), we say that the couple of measurable functions (u, v) is a weak solution of the system (4.1)–(4.3) if

$$\begin{aligned} 0 &\leq u(\mathbf{x}, t) \leq 1, \quad 0 \leq v(\mathbf{x}, t) \quad \text{for a.e. in } Q_{t_f}, \\ \xi(u) &\in L^2(0, t_f; H^1(\Omega)), \\ v &\in L^\infty(Q_{t_f}) \cap L^2(0, t_f; H^1(\Omega)), \end{aligned}$$

and for all $\varphi, \psi \in \mathcal{D}(\overline{\Omega} \times [0, t_f])$, one has

$$\begin{aligned} & - \int_{\Omega} u_0(\mathbf{x}) \varphi(\mathbf{x}, 0) d\mathbf{x} - \iint_{Q_{t_f}} u \partial_t \varphi d\mathbf{x} dt + \iint_{Q_{t_f}} \sqrt{a(u)} \Lambda(\mathbf{x}) \nabla \xi(u) \cdot \nabla \varphi d\mathbf{x} dt \\ & \quad - \iint_{Q_{t_f}} \Lambda(\mathbf{x}) \chi(u) \nabla v \cdot \nabla \varphi d\mathbf{x} dt = \iint_{Q_{t_f}} f(u) \varphi(\mathbf{x}, t) d\mathbf{x} dt, \\ & - \int_{\Omega} v_0(\mathbf{x}) \psi(\mathbf{x}, 0) d\mathbf{x} - \iint_{Q_{t_f}} v \partial_t \psi d\mathbf{x} dt + \iint_{Q_{t_f}} D(\mathbf{x}) \nabla v \cdot \nabla \psi d\mathbf{x} dt \\ & \quad = \iint_{Q_{t_f}} g(u, v) \psi d\mathbf{x} dt. \end{aligned}$$

4.3 Space-time discretization and notations

In this section, we describe the space and time discretizations of Q_{t_f} , define the approximation spaces, introduce useful properties on discrete H^1 -norms stemming from finite elements discretizations.

4.3.1 Space discretizations of Ω .

In order to discretize problem (4.1)–(4.2), we perform a finite element triangulation \mathcal{T} of the polygonal domain Ω , consisting of open bounded triangles such that $\overline{\Omega} = \bigcup_{T \in \mathcal{T}} \overline{T}$ and such that for all $T, T' \in \mathcal{T}$, $\overline{T} \cap \overline{T'}$ is either an empty set or a common vertex or edge of T and T' . We denote by \mathcal{V} the set of vertices of the discretization \mathcal{T} , located at positions $(\mathbf{x}_K)_{K \in \mathcal{V}}$, and by \mathcal{E} the set of edges of \mathcal{T} joining two vertices of \mathcal{V} , that are contained in hyperplanes of \mathbb{R}^d . The edge joining two vertices K and L is denoted by σ_{KL} .

For a given triangle $T \in \mathcal{T}$, we denote by \mathbf{x}_T the centre of gravity of T , by \mathcal{E}_T the set of the edges of T , by $h_T = \text{diam}(T)$ the diameter of T , and by ρ_T the diameter of the largest ball inscribed in the triangle T . We denote by h the size of the triangulation \mathcal{T} and by $\theta_{\mathcal{T}}$ the shape regularity assumption on the triangulation \mathcal{T} , they are given by

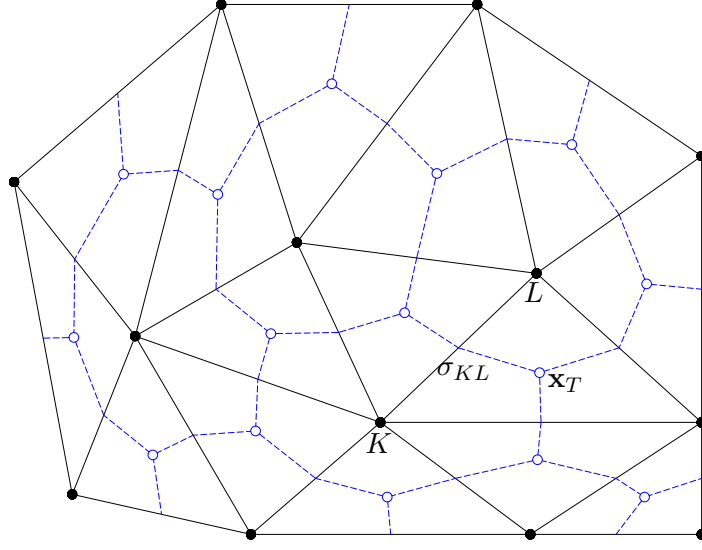
$$h_{\mathcal{T}} := \max_{T \in \mathcal{T}} h_T, \quad \theta_{\mathcal{T}} := \max_{T \in \mathcal{T}} \frac{h_T}{\rho_T}.$$

For $K \in \mathcal{V}$, we denote by \mathcal{E}_K the set of the edges having K as an extremity, and by \mathcal{T}_K the subset of \mathcal{T} including the triangles having K as a vertex. We also define a dual barycentric mesh \mathcal{M} (known as Donald dual or Median dual mesh) generated by the triangulation mesh \mathcal{T} . There is one dual element ω_K associated with each vertex $K \in \mathcal{V}$. We construct it around the vertex K by connecting the barycenter \mathbf{x}_T of each triangle $T \in \mathcal{T}_K$ with the barycenters \mathbf{x}_{σ} of the edges $\sigma \in \mathcal{E}_K$. We refer to Fig. 4.1 for an illustration of the primary and the dual barycentric mesh in a two-dimensional space. Note that $\overline{\Omega} = \bigcup_{K \in \mathcal{V}} \overline{\omega}_K$. The 2-dimensional Lebesgue measure of ω_K is denoted by m_K .

4.3.2 Discrete finite elements space $\mathcal{H}_{\mathcal{T}}$, control volumes space $\mathcal{X}_{\mathcal{M}}$.

We define two discrete functional spaces associated with each mesh of the above constructed meshes. The first one, denoted by $\mathcal{H}_{\mathcal{T}}$, is the usual \mathbb{P}_1 -finite element space corresponding to the triangular mesh \mathcal{T} , consisting of piecewise affine finite elements.

$$\mathcal{H}_{\mathcal{T}} := \{ \varphi \in C^0(\overline{\Omega}) ; \varphi|_T \in \mathbb{P}_1(\mathbb{R}), \forall T \in \mathcal{T} \} \subset H^1(\Omega).$$

FIGURE 4.1 – Triangular mesh \mathcal{T} and Donald dual mesh \mathcal{M} : dual volumes, vertices, interfaces.

The canonical basis of $\mathcal{H}_{\mathcal{T}}$ is spanned by the shape functions $(\varphi_K)_{K \in \mathcal{V}}$, such that

$$\varphi_K(\mathbf{x}_K) = 1, \quad \varphi_K(\mathbf{x}_L) = 0, \quad \text{if } L \neq K, \quad \forall K \in \mathcal{V}.$$

On the other hand, we denote by $\mathcal{X}_{\mathcal{M}}$ the discrete control volumes space consisting of piecewise constant functions on the dual mesh \mathcal{M} .

$$\mathcal{X}_{\mathcal{M}} = \{\varphi : \Omega \longrightarrow \overline{\mathbb{R}} \text{ measurable}; \varphi|_{\omega_K} \in \overline{\mathbb{R}} \text{ is constant}, \quad \forall K \in \mathcal{V}\}.$$

Given a vector $(u_K)_{K \in \mathcal{V}} \in \mathbb{R}^{\#\mathcal{V}}$ (resp. $(v_K)_{K \in \mathcal{V}} \in \mathbb{R}^{\#\mathcal{V}}$), there exists a unique function $u_{\mathcal{T}} \in \mathcal{H}_{\mathcal{T}}$ (resp. $v_{\mathcal{T}} \in \mathcal{H}_{\mathcal{T}}$) and a unique $u_{\mathcal{M}} \in \mathcal{X}_{\mathcal{M}}$ (resp. $v_{\mathcal{M}} \in \mathcal{X}_{\mathcal{M}}$) such that

$$\begin{aligned} u_{\mathcal{T}}(\mathbf{x}_K) &= u_{\mathcal{M}}(\mathbf{x}_K) = u_K, & \forall K \in \mathcal{V}, \\ v_{\mathcal{T}}(\mathbf{x}_K) &= v_{\mathcal{M}}(\mathbf{x}_K) = v_K, & \forall K \in \mathcal{V}. \end{aligned} \tag{4.6}$$

4.3.3 Time discretization of $(0, t_f)$.

For the time discretization of the interval $(0, t_f)$, we do not impose any restriction on the time step, and we restrict our study to the case of a uniform time discretization. In addition, we assume that the spatial meshes do not change with the time step. We note that all the results presented in this chapter can be extended to the case of general time discretization.

Let N be a nonnegative integer, we define the uniform time step $\Delta t = t_f / (N + 1)$, and $t_n = n\Delta t$ for all $n \in \{0, \dots, N + 1\}$, so that $t_0 = 0$, and $t_{N+1} = t_f$.

4.3.4 Space-time discretization of Q_{t_f} .

Here, we define the space and time discrete spaces $\mathcal{H}_{\mathcal{T}, \Delta t}$ and $\mathcal{X}_{\mathcal{M}, \Delta t}$ as the set of piecewise constant functions in time with values in $\mathcal{H}_{\mathcal{T}}$ and $\mathcal{X}_{\mathcal{M}}$ respectively.

$$\begin{aligned} \mathcal{H}_{\mathcal{T}, \Delta t} &= \{\varphi \in L^2(0, t_f; H^1(\Omega)), \varphi(\mathbf{x}, t) = \varphi(\mathbf{x}, t_{n+1}) \in \mathcal{H}_{\mathcal{T}}, \quad \forall t \in (t_n, t_{n+1}]\}, \\ \mathcal{X}_{\mathcal{M}, \Delta t} &= \{\varphi : Q_{t_f} \longrightarrow \overline{\mathbb{R}} \text{ measurable}, \varphi(\mathbf{x}, t) = \varphi(\mathbf{x}, t_{n+1}) \in \mathcal{X}_{\mathcal{M}}, \quad \forall t \in (t_n, t_{n+1}]\}. \end{aligned}$$

For a given $(u_K^n)_{n \in \{0, \dots, N+1\}, K \in \mathcal{V}} \in \mathbb{R}^{(N+2)\#\mathcal{V}}$ (resp. $(v_K^n)_{n \in \{0, \dots, N+1\}, K \in \mathcal{V}}$), there exists a unique function $u_{\mathcal{T}, \Delta t} \in \mathcal{H}_{\mathcal{T}, \Delta t}$ (resp. $v_{\mathcal{T}, \Delta t} \in \mathcal{H}_{\mathcal{T}, \Delta t}$) and a unique $u_{\mathcal{M}, \Delta t} \in \mathcal{X}_{\mathcal{M}, \Delta t}$ (resp. $v_{\mathcal{M}, \Delta t} \in \mathcal{X}_{\mathcal{M}, \Delta t}$) such that

$$\begin{aligned} u_{\mathcal{T}, \Delta t}(\mathbf{x}_K, t) &= u_{\mathcal{M}, \Delta t}(\mathbf{x}_K, t) = u_K^{n+1}, & \forall K \in \mathcal{V}, \forall t \in (t_n, t_{n+1}], \\ v_{\mathcal{T}, \Delta t}(\mathbf{x}_K, t) &= v_{\mathcal{M}, \Delta t}(\mathbf{x}_K, t) = v_K^{n+1}, & \forall K \in \mathcal{V}, \forall t \in (t_n, t_{n+1}]. \end{aligned} \quad (4.7)$$

4.3.5 Main property

For all $(K, L) \in \mathcal{V}^2$, we define the coefficient \mathbf{T}_{KL} by

$$\mathbf{T}_{KL} = - \int_{\Omega} \mathbf{T}(\mathbf{x}) \nabla \varphi_K(\mathbf{x}) \cdot \nabla \varphi_L(\mathbf{x}) d\mathbf{x} = \mathbf{T}_{LK}, \quad \mathbf{T}(\mathbf{x}) = \Lambda(\mathbf{x}) \text{ or } D(\mathbf{x}). \quad (4.8)$$

We have $\mathbf{T}_{KK} = - \sum_{L \neq K} \mathbf{T}_{KL}$, since $\sum_{K \in \mathcal{V}} \nabla \varphi_K = 0$. As a consequence, given $u_{\mathcal{T}}$ and $v_{\mathcal{T}}$ two elements of $\mathcal{H}_{\mathcal{T}}$, one has

$$\int_{\Omega} \mathbf{T}(\mathbf{x}) \nabla u_{\mathcal{T}} \cdot \nabla v_{\mathcal{T}} d\mathbf{x} = \sum_{\sigma_{KL} \in \mathcal{E}} a_{KL} (u_K - u_L) (v_K - v_L), \quad \mathbf{T}(\mathbf{x}) = \Lambda(\mathbf{x}) \text{ or } D(\mathbf{x}).$$

Proof. According to the definition of the approximate functions $u_{\mathcal{T}}$ and $v_{\mathcal{T}}$, one gets

$$\begin{aligned} \int_{\Omega} \mathbf{T}(\mathbf{x}) \nabla u_{\mathcal{T}} \cdot \nabla v_{\mathcal{T}} d\mathbf{x} &= \sum_{T \in \mathcal{T}} \int_T \mathbf{T}(\mathbf{x}) \nabla u_{\mathcal{T}} \cdot \nabla v_{\mathcal{T}} d\mathbf{x} \\ &= \sum_{T \in \mathcal{T}} |T| \mathbf{T}_T \left(\sum_{K \in \mathcal{V}} u_K \nabla \varphi_{K|_T} \right) \cdot \left(\sum_{L \in \mathcal{V}} v_L \nabla \varphi_{L|_T} \right) \\ &= - \sum_{K \in \mathcal{V}} \sum_{L \in \mathcal{V}, L \neq K} \mathbf{T}_{KL} u_K v_L - \sum_{K \in \mathcal{V}} \mathbf{T}_{KK} u_K v_K \\ &= - \sum_{K \in \mathcal{V}} \sum_{L \in \mathcal{V}, L \neq K} \mathbf{T}_{KL} u_K v_L + \sum_{K \in \mathcal{V}} \sum_{L \in \mathcal{V}, L \neq K} \mathbf{T}_{KL} u_K v_K \\ &= \sum_{K \in \mathcal{V}} \sum_{L \in \mathcal{V}, L \neq K} \mathbf{T}_{KL} (v_K - v_L) u_K \\ &= \sum_{\sigma_{KL} \in \mathcal{E}} \mathbf{T}_{KL} (v_K - v_L) (u_K - u_L), \end{aligned}$$

where, we have used the discrete integration by parts. \square

4.4 The nonlinear CVFE scheme

In this section, we give the discretization of system (4.1)–(4.3). The discretization defined here, is very similar to that given in chapter 3. We consider an implicit Euler scheme in time for the time evolution terms. On the other hand, the approximations in space are obtained using either the finite element approach in compliance with the Godunov scheme, or a nonclassical upwind finite volume scheme.

4.4.1 Discretization of the first equation of system (4.1)

Let us begin by introducing the discretization of the diffusive term. To do that we consider the following diffusion equation

$$\partial_t u - \operatorname{div} (\Lambda(\mathbf{x}) a(u) \nabla u) = 0. \quad (4.9)$$

By the change of variables $\vec{V} = -\nabla u$, equation (4.9) becomes

$$\partial_t u + \operatorname{div} \left(\Lambda(\mathbf{x}) a(u) \vec{V} \right) = 0,$$

which has the form of an hyperbolic PDE equation. Therefore the Godunov scheme, which was introduced in [41], may be written at the interface $\sigma_{KL} \in \mathcal{E}$ by the following expression

$$F_{KL}(u_K, u_L) = \begin{cases} \min_{u \in [u_K, u_L]} \Lambda_{KL}(u_K - u_L) a(u), & \text{if } u_K \leq u_L, \\ \max_{u \in [u_L, u_K]} \Lambda_{KL}(u_K - u_L) a(u), & \text{if } u_L \leq u_K, \end{cases}$$

where $\Lambda_{KL}(u_K - u_L)$ represents the approximation of the flux $-\Lambda \nabla u \cdot \mathbf{n}_{KL}$ at the interface σ_{KL} . F_{KL} is called the numerical flux function, it verifies the monotonicity property (see e.g. [31]). Note that, for every real number $\kappa \in \mathbb{R}$ and for every real function f , one has

$$\min(\kappa f) = \begin{cases} \kappa \min f, & \text{if } \kappa \geq 0, \\ \kappa \max f, & \text{if } \kappa < 0, \end{cases} \quad \max(\kappa f) = \begin{cases} \kappa \max f, & \text{if } \kappa \geq 0, \\ \kappa \min f, & \text{if } \kappa < 0. \end{cases}$$

As a consequence, one can conclude that the monotone flux scheme F_{KL} may be summarized by the following form

$$F_{KL}(u_K, u_L) = \begin{cases} \Lambda_{KL}(u_K - u_L) \max_{u \in [u_K, u_L]} a(u), & \text{if } \Lambda_{KL} \geq 0, \\ \Lambda_{KL}(u_K - u_L) \min_{u \in [u_K, u_L]} a(u), & \text{if } \Lambda_{KL} < 0, \\ \Lambda_{KL}(u_K - u_L) \max_{u \in [u_L, u_K]} a(u), & \text{if } \Lambda_{KL} \geq 0, \\ \Lambda_{KL}(u_K - u_L) \min_{u \in [u_L, u_K]} a(u), & \text{if } \Lambda_{KL} < 0. \end{cases}$$

Returning to the first equation of system (4.1). The discretization of the initial data u_K^0 , $K \in \mathcal{V}$ is defined by

$$u_{\mathcal{M}}^0(\mathbf{x}) = u_K^0 = \frac{1}{m_K} \int_{\omega_K} u_0(\mathbf{y}) \, d\mathbf{y}, \quad \forall \mathbf{x} \in \omega_K, \quad (4.10)$$

$$v_{\mathcal{M}}^0(\mathbf{x}) = v_K^0 = \frac{1}{m_K} \int_{\omega_K} v_0(\mathbf{y}) \, d\mathbf{y}, \quad \forall \mathbf{x} \in \omega_K, \quad (4.11)$$

and for all $K \in \mathcal{V}$, and $n \in \{0, \dots, N\}$, we define the discretization of the diffusive term by

$$\sum_{\sigma_{KL} \in \mathcal{E}_K} a_{KL}^{n+1} \Lambda_{KL}(u_K^{n+1} - u_L^{n+1}),$$

where,

$$a_{KL}^{n+1} = \begin{cases} \max_{u \in I_{KL}^{n+1}} a(u) & \text{if } \Lambda_{KL} \geq 0, \\ \min_{u \in I_{KL}^{n+1}} a(u) & \text{if } \Lambda_{KL} \leq 0, \end{cases} \quad (4.12)$$

and I_{KL}^{n+1} denotes the interval defined by

$$I_{KL}^{n+1} = \begin{cases} [u_K^{n+1}, u_L^{n+1}] & \text{if } u_K^{n+1} \leq u_L^{n+1}, \\ [u_L^{n+1}, u_K^{n+1}] & \text{otherwise.} \end{cases}$$

Let us focus on the discretization of the convective term, and recall that the function $\chi(u)$ is defined to be the product of the continuous functions $\mu(u)$ and $a(u)$. To handle the discretization of the convective term in order to obtain a robust and stable scheme, we perform a nonclassical upwind finite volume scheme which consists of considering an upwind scheme for the function $\mu(u)$ according to the discrete gradient of v , and an upwind finite volume scheme for the function $a(u)$ with respect to u . These choices of discretization are crucial to ensure the discrete maximum principle as well as the energy estimates on the approximate solutions.

We are now in a position to introduce what we call *nonlinear control volume finite element (CVFE) scheme*. For all $K \in \mathcal{V}$, and all $n \in \{0, \dots, N\}$,

$$\begin{aligned} \frac{u_K^{n+1} - u_K^n}{\Delta t} m_K + \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) \\ - \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}) = f(u_K^{n+1}) m_K, \end{aligned} \quad (4.13)$$

where the transmissibility coefficients Λ_{KL} and D_{KL} are given by equality (4.8), and μ_{KL}^{n+1} denotes an approximation of $\mu(u)$ on the interfaces of ω_K with respect to the discrete gradient of v . The term $\Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1})$ represents a numerical flux function computed at $(u_K^{n+1}, u_L^{n+1}, \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}))$. In the general case, a numerical flux function G of arguments $(a, b, c) \in \mathbb{R}^3$ is required to satisfy the following properties :

$$\begin{cases} \text{(a) } G(\cdot, b, c) \text{ is nondecreasing for all } b, c \in \mathbb{R}, \\ \text{and } G(a, \cdot, c) \text{ is nonincreasing for all } a, c \in \mathbb{R}; \\ \text{(b) } G(a, b, c) = -G(b, a, -c) \text{ for all } a, b, c \in \mathbb{R}; \\ \text{(c) } G(a, a, c) = \mu(a) c \text{ for all } a, c \in \mathbb{R}. \end{cases} \quad (4.14)$$

We give here two examples on the construction of μ_{KL}^{n+1} such that the numerical flux function satisfies properties (4.14). The first example consists of taking the Engquist-Osher scheme [65] and the second example consists of taking the Godunov scheme (see e.g. [32, 40]).

$$\bullet \mu_{KL}^{n+1} = \begin{cases} \mu_{\downarrow}(u_K^{n+1}) + \mu_{\uparrow}(u_L^{n+1}), & \text{if } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) \geq 0, \\ \mu_{\uparrow}(u_K^{n+1}) + \mu_{\downarrow}(u_L^{n+1}), & \text{if } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) < 0. \end{cases}$$

The functions μ_{\uparrow} and μ_{\downarrow} are given by

$$\mu_{\uparrow}(z) := \int_0^z (\mu'(s))^+ ds, \quad \mu_{\downarrow}(z) := - \int_0^z (\mu'(s))^- ds.$$

Herein, $s^+ = \max(s, 0)$ and $s^- = \max(-s, 0)$.

$$\bullet \mu_{KL}^{n+1} = \begin{cases} \max_{[u_K^{n+1}, u_L^{n+1}]} \mu(u), & \text{if } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) \geq 0, \\ \min_{[u_L^{n+1}, u_K^{n+1}]} \mu(u), & \text{if } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) \geq 0, \\ \max_{[u_L^{n+1}, u_K^{n+1}]} \mu(u), & \text{if } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) < 0, \\ \min_{[u_K^{n+1}, u_L^{n+1}]} \mu(u), & \text{if } \Lambda_{KL} (v_K^{n+1} - v_L^{n+1}) < 0. \end{cases}$$

Remark 3. In view of properties (4.14) of the numerical flux function, one can deduce that the approximate numerical flux function $\mu_{KL}^{n+1} = z(u_K^{n+1}, u_L^{n+1})$ must be nonincreasing (resp. non-decreasing) with respect to u_K^{n+1} and nondecreasing (resp. nonincreasing) with respect to u_L^{n+1} when $\Lambda_{KL}(v_K^{n+1} - v_L^{n+1}) \geq 0$ (resp. $\Lambda_{KL}(v_K^{n+1} - v_L^{n+1}) \leq 0$). We also have, using property (4.14)(c), that $\mu_{KK}^{n+1} = z(u_K^{n+1}, u_K^{n+1}) = \mu(u_K^{n+1})$.

4.4.2 Discretization of the second equation of system (4.1)

We now focus on the discretization of the second equation of system (4.1). We note that a classical discretization of this equation is given by the following form

$$m_K \frac{v_K^{n+1} - v_K^n}{\Delta t} + \sum_{\sigma_{KL} \in \mathcal{E}_K} D_{KL} (v_K^{n+1} - v_L^{n+1}) = m_K (\alpha u_K^n - \beta v_K^{n+1}). \quad (4.15)$$

However, we have seen in chapter 3 that this discretization does not guaranty the discrete maximum principle without any restriction on the transmissibility coefficients, for instance, one can get the discrete maximum principle by assuming that all the transmissibility coefficients D_{KL} are nonnegative (e.g. see [50]).

Here, we propose a numerical discretization in order to ensure the discrete maximum principle without any restriction on the transmissibility coefficients. To do this, we introduce the following set of functions : $\eta(v)$, $p(v)$, $\Gamma(v)$ and $\phi(v)$ defined by

$$\eta(v) = \max(0, \min(v, 1)), \quad (4.16)$$

$$p(v) = \int_1^v \frac{1}{\eta(s)} ds = \begin{cases} \ln(v) & \text{if } v \in (0, 1), \\ v - 1 & \text{if } v \geq 1, \end{cases} \quad (4.17)$$

$$\Gamma(v) = \int_1^v p(s) ds = \begin{cases} v \ln(v) - v + 1 & \text{if } v \in [0, 1), \\ \frac{(v-1)^2}{2} & \text{if } v \geq 1, \end{cases} \quad (4.18)$$

$$\phi(v) = \int_0^v \frac{1}{\sqrt{\eta(s)}} ds = \begin{cases} \frac{\sqrt{v}-1}{2} & \text{if } v \in [0, 1), \\ v - 1 & \text{if } v \geq 1. \end{cases} \quad (4.19)$$

In the sequel, we adopt the convention

$$\eta(v)p(v) = 0 \quad \text{if } v \leq 0. \quad (4.20)$$

We mention the discretization of the second equation of (4.1); specifically, we have

$$m_K \frac{v_K^{n+1} - v_K^n}{\Delta t} - \sum_{\sigma_{KL} \in \mathcal{E}_K} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1})) = m_K (\alpha u_K^n - \beta v_K^{n+1}), \quad (4.21)$$

where, denoting by $J_{KL}^{n+1} = [\min(v_K^{n+1}, v_L^{n+1}), \max(v_K^{n+1}, v_L^{n+1})]$, we have set

$$\eta_{KL}^{n+1} = \begin{cases} \max_{s \in J_{KL}^{n+1}} \eta(s) & \text{if } D_{KL} \geq 0, \\ \min_{s \in J_{KL}^{n+1}} \eta(s) & \text{if } D_{KL} < 0. \end{cases}$$

Taking into account properties (4.14) of the numerical flux function G , one can deduce that this scheme, whose construction is based on finite elements for the diffusive term and an upstream finite volume for the convective term, can be interpreted as a finite volume scheme. Indeed, denoting by

$$\begin{aligned} F_{KL}^{n+1} &= \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) - \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}), \\ \Phi_{KL}^{n+1} &= D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1})). \end{aligned}$$

We have $a_{KL}^{n+1} = a_{LK}^{n+1}$, $\eta_{KL}^{n+1} = \eta_{LK}^{n+1}$, and $\mu_{KL}^{n+1} = \mu_{LK}^{n+1}$. Thus, we obtain the locally conservative form

$$\begin{cases} F_{KL}^{n+1} + F_{LK}^{n+1} = 0 = \Phi_{KL}^{n+1} + \Phi_{LK}^{n+1}, & \text{for all } \sigma_{KL} \in \mathcal{E}, \\ m_K \frac{u_K^{n+1} - u_K^n}{\Delta t} + \sum_{\sigma_{KL} \in \mathcal{E}_K} F_{KL}^{n+1} = f(u_K^{n+1}) m_K, & \text{for all } K \in \mathcal{V}, \\ m_K \frac{v_K^{n+1} - v_K^n}{\Delta t} + \sum_{\sigma_{KL} \in \mathcal{E}_K} \Phi_{KL}^{n+1} = g(u_K^n, v_K^{n+1}) m_K, & \text{for all } K \in \mathcal{V}. \end{cases}$$

4.4.3 Main result

Let $(\mathcal{T}_m)_{m \geq 1}$ be a sequence of triangulations of Ω such that

$$h_m = \max_{T \in \mathcal{T}_m} \text{diam}(T) \rightarrow 0 \text{ as } m \rightarrow \infty.$$

We assume that there exists a constant $\theta > 0$ such that

$$\theta_{\mathcal{T}_m} \leq \theta, \quad \forall m \geq 1.$$

As before, a sequence of dual meshes $(\mathcal{M}_m)_{m \geq 1}$ is given.

Let $(N_m)_m$ be an increasing sequence of integers, then we define the corresponding sequence of time steps $(\Delta t_m)_m$ such that $\Delta t_m \rightarrow 0$ as $m \rightarrow \infty$. The intention of this chapter is to prove the following main result.

Theorem 4.2. *Let $(u_{\mathcal{M}_m, \Delta t_m}, v_{\mathcal{M}_m, \Delta t_m})_m$ be a sequence of solutions to the scheme (4.13)–(4.21), such that $0 \leq u_{\mathcal{M}_m, \Delta t_m} \leq 1$ and $0 \leq v_{\mathcal{M}_m, \Delta t_m}$ for almost everywhere in Q_{t_f} , then*

$$u_{\mathcal{M}_m, \Delta t_m} \rightarrow u \text{ and } v_{\mathcal{M}_m, \Delta t_m} \rightarrow v \quad \text{a.e. in } Q_{t_f} \text{ as } m \rightarrow \infty,$$

where the couple (u, v) is a weak solution to the system (4.1)–(4.3) in the sense of Definition 4.1.

4.5 Discrete properties, a priori estimates and existence

In this section, we first bring up some preliminary lemmas presented by Cances and Guichard [16], that we reproduce them here for clarity. We then establish the discrete maximum principle which is the basis to the analysis that we are going to perform. Next, we carry out the *a priori* estimates necessary to prove the existence of a discrete solution to the discrete problem (4.13)–(4.21). The maximum principle and the *a priori* estimates are crucial to prove the convergence of the scheme (4.13)–(4.21) towards the weak solution.

Lemma 4.3. *Let $(u_K^{n+1})_{K,n} \in \mathbb{R}^{(N+1)\#\mathcal{V}}$ (resp. $(v_K^{n+1})_{K,n} \in \mathbb{R}^{(N+1)\#\mathcal{V}}$), then denoting by $\xi_{\mathcal{T}, \Delta t}$ (resp. $\phi_{\mathcal{T}, \Delta t}$) the unique function of $\mathcal{H}_{\mathcal{T}, \Delta t}$ with nodal values $(\xi(u_K^{n+1}))_{K,n} \in \mathbb{R}^{(N+1)\#\mathcal{V}}$ (resp. $(\phi(v_K^{n+1}))_{K,n} \in \mathbb{R}^{(N+1)\#\mathcal{V}}$), one has*

$$\begin{aligned} & \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 \\ & \geq \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} (\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2 = \iint_{Q_{t_f}} \Lambda \nabla \xi_{\mathcal{T}, \Delta t} \cdot \nabla \xi_{\mathcal{T}, \Delta t} dx dt. \end{aligned} \tag{4.22}$$

and

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2 \\ \geq \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} (\phi(v_K^{n+1}) - \phi(v_L^{n+1}))^2 = \iint_{Q_{t_f}} D \nabla \phi_{\mathcal{T}, \Delta t} \cdot \nabla \phi_{\mathcal{T}, \Delta t} dx dt. \end{aligned} \quad (4.23)$$

Proof. Using the mean value theorem since ξ is a differentiable function, one can deduce the existence of $c \in I_{KL}^{n+1}$, such that

$$(\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2 = a(c) (u_K^{n+1} - u_L^{n+1})^2 \quad \forall \sigma_{KL} \in \mathcal{E}.$$

Multiplying this equality by Λ_{KL} , one gets

$$\Lambda_{KL} (\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2 \leq \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2,$$

for the case where $\Lambda_{KL} \geq 0$ since a_{KL}^{n+1} is the maximum of $a(u)$ on the interval I_{KL}^{n+1} . On the other hand, when $\Lambda_{KL} \leq 0$, one has that a_{KL}^{n+1} is the minimum of $a(u)$ on the interval I_{KL}^{n+1} , thus $\Lambda_{KL} a(c) \leq \Lambda_{KL} a_{KL}^{n+1}$ and property (4.22) holds. The proof of inequality (4.23) is similar. This ends the proof of the lemma. \square

Let $T \in \mathcal{T}$, and let $(K, L) \in \mathcal{V}^2$, we denote by

$$\lambda_{KL}^T := - \int_T \Lambda(\mathbf{x}) \nabla \varphi_K(\mathbf{x}) \cdot \nabla \varphi_L(\mathbf{x}) d\mathbf{x} = \lambda_{L,K}^T.$$

In particular, one has $\Lambda_{KL} = - \sum_{T \in \mathcal{T}} \int_T \Lambda \nabla \varphi_K \cdot \nabla \varphi_L d\mathbf{x} = \sum_{T \in \mathcal{T}} \lambda_{KL}^T$, for all $\sigma_{KL} \in \mathcal{E}$.

Lemma 4.4. *Let $\Psi_{\mathcal{T}} = \sum_{K \in \mathcal{V}} \psi_K \varphi_K \in \mathcal{H}_{\mathcal{T}}$, then there exists a quantity C_0 depending only on Λ , $\theta_{\mathcal{T}}$ such that*

$$\sum_{\sigma_{KL} \in \mathcal{E}} \sum_{T \in \mathcal{T}} |\lambda_{KL}^T| (\psi_K - \psi_L)^2 \leq C_0 \int_{\Omega} \Lambda \nabla \Psi_{\mathcal{T}} \cdot \nabla \Psi_{\mathcal{T}} d\mathbf{x}. \quad (4.24)$$

Proof. In the proof below, unless specified, C denotes a generic quantity depending only on Λ and $\theta_{\mathcal{T}}$.

In order to prove inequality (4.24), it suffices to prove that

$$\sum_{\sigma_{KL} \in \mathcal{E}} \sum_{T \in \mathcal{T}} |\lambda_{KL}^T| (\psi_K - \psi_L)^2 \leq C \|\nabla \Psi_{\mathcal{T}}\|_{(L^2(\Omega))^2}^2,$$

since we have, using the assumption on Λ , that

$$\|\nabla \Psi_{\mathcal{T}}\|_{L^2(\Omega)}^2 = \int_{\Omega} \nabla \Psi_{\mathcal{T}} \cdot \nabla \Psi_{\mathcal{T}} d\mathbf{x} \leq \frac{1}{\Lambda_-} \int_{\Omega} \Lambda \nabla \Psi_{\mathcal{T}} \cdot \nabla \Psi_{\mathcal{T}} d\mathbf{x}.$$

Thanks to the Cauchy-Schwarz inequality, one has

$$|\lambda_{KL}^T| \leq \Lambda_- \|\nabla \varphi_K\|_{(L^2(T))^2} \|\nabla \varphi_L\|_{(L^2(T))^2}. \quad (4.25)$$

Now, using a classical inequality stemming from finite element properties (see e.g. [27, 11]), one gets

$$\|\nabla \varphi_K\|_{L^2(T)}^2 \leq c\theta_T \frac{|T|}{(h_T)^2}, \quad \forall K \in \mathcal{V}, \forall T \in \mathcal{T}, \quad (4.26)$$

where c is an absolute constant.

Plugging inequality (4.26) into (4.25), one can conclude that

$$|\lambda_{KL}^T| \leq C \frac{|T|}{(h_T)^2}, \quad \forall T \in \mathcal{T}, \forall \sigma_{KL} \in \mathcal{E}_T. \quad (4.27)$$

This implies that

$$\sum_{\sigma_{KL} \in \mathcal{E}_T} |\lambda_{KL}^T| (\psi_K - \psi_L)^2 \leq C \frac{|T|}{(h_T)^2} \sum_{\sigma_{KL} \in \mathcal{E}_T} (\psi_K - \psi_L)^2.$$

Finally, it follows from the analysis carried out for example in [11] that for all $T \in \mathcal{T}$ with K, L , and M being its vertices

$$\frac{|T|}{(h_T)^2} \left((\psi_K - \psi_L)^2 + (\psi_K - \psi_M)^2 + (\psi_L - \psi_M)^2 \right) \leq C \|\nabla \Psi_T\|_{(L^2(T))^2}^2.$$

Since $\sigma_{KL} \in \mathcal{E}$ is shared by at most two triangles, one has

$$\sum_{T \in \mathcal{T}} \sum_{\sigma_{KL} \in \mathcal{E}_T} |\lambda_{KL}^T| (\psi_K - \psi_L)^2 \leq C \sum_{T \in \mathcal{T}} \|\nabla \Psi_T\|_{L^2(T)}^2 = C \|\nabla \Psi_T\|_{(L^2(\Omega))^2}^2,$$

this concludes the proof of the lemma. \square

Lemma 4.5. *There exists a quantity C_1 depending only on Λ and θ_T such that*

$$\sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |\Lambda_{KL}| a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 \leq C_1 \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2. \quad (4.28)$$

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |D_{KL}| \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2 \\ \leq C_1 \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2. \end{aligned} \quad (4.29)$$

Proof. We denote by $\mathcal{E}^- := \{\sigma_{KL} \in \mathcal{E}; \Lambda_{KL} < 0\}$, then since $|\mathbf{x}| = \mathbf{x} + 2\mathbf{x}^-, \mathbf{x}^- = \max(-\mathbf{x}, 0)$, one has

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |\Lambda_{KL}| a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 &= \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 \\ &\quad + 2 \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}^-} |\Lambda_{KL}| a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2. \end{aligned}$$

Now, from the definition (4.8) of a_{KL}^{n+1} , there exists $c \in I_{KL}^{n+1}$ such that

$$(\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2 = a(c)(u_K^{n+1} - u_L^{n+1})^2 \geq a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2, \quad \forall \sigma_{KL} \in \mathcal{E}^-.$$

Therefore, it yields

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |\Lambda_{KL}| a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 &\leq \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 \\ &\quad + 2 \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |\Lambda_{KL}| (\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2. \end{aligned} \quad (4.30)$$

Lemma 4.4 ensures the existence of a quantity $C_0 > 0$ ($= C_0(\Lambda, \theta_T)$) such that

$$\begin{aligned} \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |\Lambda_{KL}| (\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2 \\ \leq \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \sum_{T \in \mathcal{T}} |\lambda_{KL}| (\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2 \leq C_0 \int_{Q_T} \Lambda \nabla \xi_{T, \Delta t} \cdot \nabla \xi_{T, \Delta t} dx dt, \end{aligned}$$

and from lemma 4.3, we deduce that

$$\sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |\Lambda_{KL}| (\xi(u_K^{n+1}) - \xi(u_L^{n+1}))^2 \leq C_0 \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2. \quad (4.31)$$

Plugging estimate (4.31) into estimate (4.30), then estimate (4.28) holds with $C_1 = 1 + 2C_0$. The proof of (4.29) is similar. This ends the proof of the lemma. \square

4.5.1 Discrete maximum principle

Lemma 4.6. *Let $(u_K^{n+1})_{K \in \mathcal{V}, n \in \{0, \dots, N\}}$ be a solution to the CVFE scheme (4.13). Then, for all $K \in \mathcal{V}_h$, and all $n \in \{0, \dots, N+1\}$, we have $0 \leq u_K^n \leq 1$ and $v_K^n \geq 0$.*

Proof. We want to show this property using the induction on n . The property is true for $n = 0$ thanks to the definition (4.10) of u_K^0 and to the assumption on u_0 . Now, assume that the claim is true up to time step n . Consider a dual control volume ω_K such that $u_K^{n+1} = \min_{L \in \mathcal{V}} \{u_L^{n+1}\}$, we want

to show that $u_K^{n+1} \geq 0$ i.e. $(u_K^{n+1})^- = 0$.

Multiplying the scheme (4.13) by $-(u_K^{n+1})^-$, one has

$$\begin{aligned} -m_K \frac{u_K^{n+1} - u_K^n}{\Delta t} (u_K^{n+1})^- - \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) (u_K^{n+1})^- \\ + \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}) (u_K^{n+1})^- = -m_K f(u_K^{n+1}) (u_K^{n+1})^-, \end{aligned} \quad (4.32)$$

for which, we use the extension by $f(0) \geq 0$ of the function f for $u \leq 0$, then the right hand side of equation (4.32) is less or equal to zero.

In view of definition (4.12) of a_{KL}^{n+1} , and of the fact that $a(u) = 0$ for every $u \leq 0$, one has

$$a_{KL}^{n+1} (u_K^{n+1})^- = 0, \quad \text{if } \Lambda_{KL} \leq 0.$$

Indeed, assume that $u_K^{n+1} \geq 0$ then $(u_K^{n+1})^- = 0$ and $a_{KL}^{n+1} (u_K^{n+1})^- = 0$. Nevertheless, assume that $u_K^{n+1} < 0$ then $a_{KL}^{n+1} = \min_{u \in I_{KL}^{n+1}} a(u) = 0$ if $\Lambda_{KL} \leq 0$.

As a consequence, the second term in the left hand side of equation (4.32) reads to

$$- \sum_{\sigma_{KL} \in \mathcal{E}_K} a_{KL}^{n+1} (\Lambda_{KL})^+ (u_K^{n+1} - u_L^{n+1}) (u_K^{n+1})^- \geq 0,$$

Let us now focus on the third term of equation (4.32), and denote by \mathcal{A} this term. Using the fact that $a_{KL}^{n+1} (u_K^{n+1})^- = 0$ for $\Lambda_{KL} \leq 0$, then \mathcal{A} rewrites

$$\begin{aligned} \mathcal{A} = & \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL}^+ \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1})^+ (u_K^{n+1})^- \\ & - \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL}^+ \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1})^- (u_K^{n+1})^-. \end{aligned}$$

The second term of \mathcal{A} is nonpositive, but in view of Remark 3 on the approximation μ_{KL}^{n+1} and in view of the extension by zero of the function μ for $u \leq 0$ (since $\mu(0) = 0$), one can deduce that

$$\mu_{KL}^{n+1} \Lambda_{KL}^+ (v_K^{n+1} - v_L^{n+1})^- (u_K^{n+1})^- \leq \mu(u_K^{n+1}) \Lambda_{KL}^+ (v_K^{n+1} - v_L^{n+1})^- (u_K^{n+1})^- = 0,$$

thus, the second term of \mathcal{A} is equal to zero and consequently $\mathcal{A} \geq 0$ since the first term of \mathcal{A} is nonnegative.

Finally, we use the identity $u_K^{n+1} = (u_K^{n+1})^+ - (u_K^{n+1})^-$ and the nonnegativity of u_K^n , one can deduce from equation (4.32) that $(u_K^{n+1})^- = 0$. According to the choice of the dual control volume ω_K , then $\min_{L \in \mathcal{V}} \{u_L^{n+1}\}$ is non-negative. Consequently,

$$u_K^n \geq 0, \quad \forall K \in \mathcal{V}, \text{ and all } n \in \{0, \dots, N+1\}.$$

In order to prove by induction that $u_K^n \leq 1$, $\forall K \in \mathcal{V}$, $\forall n \in \{0, \dots, N+1\}$, we proceed in the same way as before, so that we assume that the claim $u_K^n \leq 1$ is true for all $K \in \mathcal{V}$ and we consider a dual control volume ω_K such that $u_K^{n+1} = \max_{L \in \mathcal{V}} \{u_L^{n+1}\}$, we want to show that $u_K^{n+1} \leq 1$.

For the mentioned claim, we multiply equation (4.13) by $(u_K^{n+1} - 1)^+$, one gets

$$\begin{aligned} m_K \frac{u_K^{n+1} - u_K^n}{\Delta t} (u_K^{n+1} - 1)^+ + \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) (u_K^{n+1} - 1)^+ \\ - \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}) (u_K^{n+1} - 1)^+ = 0, \end{aligned} \quad (4.33)$$

for which, we have use the extension by $f(1) \leq 0$ of the function f for $u \geq 1$.

In view of definition (4.12) of a_{KL}^{n+1} , and of the fact that $a(u) = 0$ for every $u \geq 1$, one has

$$a_{KL}^{n+1} (u_K^{n+1} - 1)^+ = 0, \quad \text{if } \Lambda_{KL} \leq 0.$$

the second term in the left hand side of equation (4.33) reads to

$$\sum_{\sigma_{KL} \in \mathcal{E}_K} a_{KL}^{n+1} (\Lambda_{KL})^+ (u_K^{n+1} - u_L^{n+1}) (u_K^{n+1} - 1)^+ \geq 0. \quad (4.34)$$

Using estimate (4.34) and denoting by \mathcal{A}_1 the third term of equation (4.33), one has

$$\mathcal{A}_1 \geq - \sum_{\sigma_{KL} \in \mathcal{E}_K} (\Lambda_{KL})^+ \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1})^+ (u_K^{n+1} - 1)^+. \quad (4.35)$$

The right hand side of inequality (4.35) is nonpositive, but in view of Remark 3 on the approximation μ_{KL}^{n+1} and in view of the extension by zero of the function μ for $u \geq 1$ (since $\mu(1) = 0$), one can deduce that

$$\mu_{KL}^{n+1} \Lambda_{KL}^+ (v_K^{n+1} - v_L^{n+1})^+ (u_K^{n+1} - 1)^+ \leq \mu(u_K^{n+1}) \Lambda_{KL}^+ (v_K^{n+1} - v_L^{n+1})^+ (u_K^{n+1} - 1)^+ = 0,$$

and consequently, one gets that \mathcal{A}_1 is nonnegative.

Using the identity $(u_K^{n+1} - 1) = (u_K^{n+1} - 1)^+ - (u_K^{n+1} - 1)^-$ and that $u_K^n \leq 1$, one can deduce from equation (4.33) that $(u_K^{n+1} - 1)^+ = 0$. According to the choice of the dual control volume ω_K , then $\max_{L \in \mathcal{V}} \{u_L^{n+1}\}$ is non-negative. Consequently,

$$u_K^n \leq 1, \quad \forall K \in \mathcal{V}, \text{ and all } n \in \{0, \dots, N+1\}.$$

Let us prove now that $v_K^{n+1} \geq 0$ for all $K \in \mathcal{V}$. Let $K_m \in \mathcal{V}$ be such that $v_{K_m}^{n+1} = \min_{K \in \mathcal{V}} v_K^{n+1}$, and assume that $v_{K_m}^{n+1} \leq 0$. Thanks to the convention (4.20), we can claim that

$$D_{K_m L} \eta_{K_m L}^{n+1} (p(v_{K_m}^{n+1}) - p(v_L^{n+1})) = 0 \quad \text{if } D_{K_m L} \leq 0,$$

and that

$$D_{K_m L} \eta_{K_m L}^{n+1} (p(v_{K_m}^{n+1}) - p(v_L^{n+1})) \leq 0 \quad \text{if } D_{K_m L} \geq 0.$$

Therefore, the scheme (4.21) together with the positivity of $u_{K_m}^n$ yields

$$v_{K_m}^{n+1} \geq \frac{v_{K_m}^n}{1 + \beta \Delta t} \geq 0.$$

This achieves the proof of Lemma 4.6. □

4.5.2 Entropy estimates on $v_{\mathcal{M}, \Delta t}$

Lemma 4.7. *There exists $C > 0$ depending only on $\|v_0\|_{L^2(\Omega)}$, Ω , t_f , α and β such that, for all $n^* \in \{0, \dots, N\}$, one has*

$$\sum_{K \in \mathcal{V}} m_K \Gamma(v_K^{n^*+1}) + \sum_{n=0}^{n^*} \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2 \leq C.$$

Proof. It follows from Jensen's inequality – recall that Γ is convex – that

$$\sum_{K \in \mathcal{V}} m_K \Gamma(v_K^0) \leq \int_{\Omega} \Gamma(v_0(\mathbf{x})) d\mathbf{x}.$$

Since $\Gamma(v) \leq (v - 1)^2$ for all $v \geq 0$, we obtain that

$$\sum_{K \in \mathcal{V}} m_K \Gamma(v_K^0) \leq \int_{\Omega} (v_0(\mathbf{x}) - 1)^2 d\mathbf{x} \leq C. \quad (4.36)$$

Multiplying the scheme (4.21) by $p(v_K^{n+1}) \Delta t$ and summing of $K \in \mathcal{V}$ and $n = 0, \dots, n^*$ provides

$$\mathcal{A} + \mathcal{B} = \mathcal{C}, \quad (4.37)$$

where we have set

$$\begin{aligned}\mathcal{A} &= \sum_{n=0}^{n^*} \sum_{K \in \mathcal{V}} m_K (v_K^{n+1} - v_K^n) p(v_K^{n+1}), \\ \mathcal{B} &= \sum_{n=0}^{n^*} \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2, \\ \mathcal{C} &= \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} m_K (\alpha u_K^n - \beta v_K^{n+1}) p(v_K^{n+1}).\end{aligned}$$

Since, thanks to Lemma 4.6, u_K^n is nonnegative for all $K \in \mathcal{V}$ and all $n \geq 0$, and since $p(v) \leq (v - 1)$ for all $v \geq 0$ (with the convention $p(0) = -\infty$), one has

$$\alpha u_K^n p(v_K^{n+1}) \leq \alpha u_K^n (v_K^{n+1} - 1).$$

On the other hand, there exists an absolute constant c^* such that $vp(v) \geq (v - 1)^2 - c^*$ for all $v \geq 0$. Therefore,

$$\beta v_K^{n+1} p(v_K^{n+1}) \geq \beta (v_K^{n+1} - 1)^2 - c^*.$$

As a consequence, we obtain that

$$\mathcal{C} \leq t_f |\Omega| c^* + \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} m_K (\alpha u_K^n (v_K^{n+1} - 1) - \beta (v_K^{n+1} - 1)^2).$$

Using the weighted Young's inequality $\alpha ab \leq \beta b^2 + \frac{\alpha^2}{4\beta} a^2$ for all $(a, b) \in \mathbb{R}^2$ provides

$$\alpha u_K^n (v_K^{n+1} - 1) - \beta (v_K^{n+1} - 1)^2 \leq \frac{\alpha^2}{4\beta} u_K^n \leq \frac{\alpha^2}{4\beta}$$

thanks to Lemma 4.6. Hence, we obtain that

$$\mathcal{C} \leq t_f |\Omega| \left(c^* + \frac{\alpha^2}{4\beta} \right). \quad (4.38)$$

The function p being increasing, an elementary convexity inequality provides that

$$(a - b)p(a) \geq \Gamma(a) - \Gamma(b), \quad \forall (a, b) \in (\mathbb{R}_+)^2,$$

ensuring that

$$\mathcal{A} \geq \sum_{n=0}^{n^*} \sum_{K \in \mathcal{V}} m_K (\Gamma(v_K^{n+1}) - \Gamma(v_K^n)) = \sum_{K \in \mathcal{V}} m_K (\Gamma(v_K^{n^*+1}) - \Gamma(v_K^0)). \quad (4.39)$$

Using (4.38), (4.39) and (4.36) in (4.37) concludes the proof of Lemma 4.7. \square

Lemma 4.8. *There exists C depending only on $\|v_0\|_{L^2(\Omega)}$, Ω , t_f , α , β , D_{\pm} , Λ_+ and $\theta_{\mathcal{T}}$, such that*

$$\iint_{Q_{t_f}} \Lambda(\mathbf{x}) \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) \cdot \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) d\mathbf{x} dt = \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} (v_K^{n+1} - v_L^{n+1})^2 \leq C. \quad (4.40)$$

Proof. Thanks to Lemmas 4.3 and 4.7, we know that

$$\sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} (\phi(v_K^{n+1}) - \phi(v_L^{n+1}))^2 \leq \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2 \leq C.$$

Therefore, it follows from Lemma 4.4 that

$$\sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |D_{KL}| (\phi(v_K^{n+1}) - \phi(v_L^{n+1}))^2 \leq C.$$

since $(\phi(v_K^{n+1}) - \phi(v_L^{n+1}))^2 \geq (v_K^{n+1} - v_L^{n+1})^2$, we obtain that

$$\begin{aligned} \iint_{Q_{t_f}} D(\mathbf{x}) \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) \cdot \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) d\mathbf{x} dt &= \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} (v_K^{n+1} - v_L^{n+1})^2 \\ &\leq \sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |D_{KL}| (v_K^{n+1} - v_L^{n+1})^2 \leq C. \end{aligned}$$

It only remains to check that

$$\iint_{Q_{t_f}} \Lambda(\mathbf{x}) \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) \cdot \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) d\mathbf{x} dt \leq \frac{\Lambda_+}{D_-} \iint_{Q_{t_f}} D(\mathbf{x}) \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) \cdot \nabla v_{\mathcal{T}, \Delta t}(\mathbf{x}, t) d\mathbf{x} dt$$

in order to conclude the proof of Lemma 4.8. \square

Remark 4. A careful analysis allows to prove, after a more involving proof very similar to the one of Lemma 4.7, that the constant C depends neither on D_+ nor on $\theta_{\mathcal{T}}$. The cornerstone of the corresponding proof is the inequality

$$D_{KL} \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1})) (v_K^{n+1} - v_L^{n+1}) \geq D_{KL} (v_K^{n+1} - v_L^{n+1})^2$$

that relies on the definition of η_{KL}^{n+1} and on the link between the functions η and p .

4.5.3 Energy estimates on $u_{\mathcal{M}, \Delta t}$

In the following, C denotes a "generic" constant, which need not have the same value throughout the proofs.

Proposition 4.9. *Let $(u_K^{n+1}, v_K^{n+1})_{K \in \mathcal{V}, n \in \{0, \dots, N\}}$ be a solution to the scheme (4.13)-(4.21). There exists a constant $C > 0$ depending only on $\|v_0\|_{L^2(\Omega)}$, Ω , t_f , α , β , Λ , D , and $\theta_{\mathcal{T}}$ such that*

$$\sum_{K \in \mathcal{V}} m_K \left(u_K^{n^*+1} \right)^2 + \sum_{n=0}^{n^*} \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 \leq C. \quad (4.41)$$

Proof. We multiply equation (4.13) by $\Delta t u_K^{n+1}$ and sum over $K \in \mathcal{V}$ and $n \in \{0, \dots, n^*\}$. This yields

$$E_1 + E_2 + E_3 = E_4, \quad (4.42)$$

where

$$\begin{aligned}
E_1 &= \sum_{n=0}^{n^*} \sum_{K \in \mathcal{V}} m_K (u_K^{n+1} - u_K^n) u_K^{n+1}, \\
E_2 &= \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) u_K^{n+1}, \\
E_3 &= - \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}) u_K^{n+1}, \\
E_4 &= \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} m_K f(u_K^{n+1}) u_K^{n+1}.
\end{aligned}$$

For the time evolution term, we use the following inequality : $(a - b) a \geq \frac{1}{2} (a^2 - b^2)$, $\forall a, b \in \mathbb{R}$, to get

$$E_1 \geq \frac{1}{2} \sum_{n=0}^{n^*} \sum_{K \in \mathcal{V}_h} m_K \left((u_K^{n+1})^2 - (u_K^n)^2 \right) = \frac{1}{2} \sum_{K \in \mathcal{V}_h} m_K \left((u_K^{n^*+1})^2 - (u_K^0)^2 \right). \quad (4.43)$$

Next, for the diffusion term, we reorganize the sum over the edges, we find

$$\begin{aligned}
E_2 &= \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) u_K^{n+1} \\
&= \sum_{n=0}^{n^*} \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2.
\end{aligned} \quad (4.44)$$

Similarly, we reorganize the sum over the edges for the convection term, we obtain

$$\begin{aligned}
E_3 &= - \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}) u_K^{n+1} \\
&= - \sum_{n=0}^{n^*} \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}) (u_K^{n+1} - u_L^{n+1}).
\end{aligned}$$

Using the weighted Young inequality, we deduce

$$\begin{aligned}
|E_3| &\leq C \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} |\Lambda_{KL}| a_{KL}^{n+1} |v_K^{n+1} - v_L^{n+1}| |u_K^{n+1} - u_L^{n+1}| \\
&\leq C \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} |\Lambda_{KL}| (v_K^{n+1} - v_L^{n+1})^2 \\
&\quad + \frac{1}{2C_1} \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} |\Lambda_{KL}| a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2.
\end{aligned}$$

Thanks to estimates (4.24) and (4.28), one has

$$\begin{aligned} |E_3| &\leq C \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} (v_K^{n+1} - v_L^{n+1})^2 \\ &\quad + \frac{1}{2} \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2. \end{aligned}$$

Therefore, Lemma 4.8 provides that

$$|E_3| \leq C + \frac{1}{2} \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}} \sum_{\sigma_{KL} \in \mathcal{E}_K} \Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2. \quad (4.45)$$

Finally, for the reaction term, since $0 \leq u_K^{n+1} \leq 1$ thanks to Lemma 4.6, one has

$$E_4 = \sum_{n=0}^{n^*} \Delta t \sum_{K \in \mathcal{V}_h} m_K f(u_K^{n+1}) u_K^{n+1} \leq |\Omega| \|f\|_{L^\infty(0,1)} t_f. \quad (4.46)$$

Plugging estimates (4.43)–(4.46) into equation (4.42), one can deduce that the estimate (4.41) holds. \square

4.5.4 Enhanced estimate on $v_{\mathcal{M}, \Delta t}$

The goal of this section is to prove a refined estimate on $v_{\mathcal{M}, \Delta t}$ inspired from [16, Lemma 3.10], claiming that either $v_{\mathcal{M}, \Delta t}$ is constant equal to 0, or $v_{\mathcal{M}, \Delta t} \geq r_h > 0$ for some r_h depending on the discretization parameters. The first step consists in bounding from below the $L^\infty((0, t_f); L^1(\Omega))$ norm of $v_{\mathcal{M}, \Delta t}$.

Lemma 4.10. *Assume that $\int_\Omega u_0(\mathbf{x}) d\mathbf{x} > 0$ or $\int_\Omega v_0(\mathbf{x}) d\mathbf{x} > 0$, then there exists $\kappa > 0$ depending on the discretization and on the data such that*

$$\int_\Omega v_{\mathcal{M}, \Delta t}(\mathbf{x}, t) d\mathbf{x} \geq \kappa, \quad \forall t \in [0, t_f].$$

Proof. Summing equation (4.21) over $K \in \mathcal{V}$ ensures that

$$\sum_{K \in \mathcal{V}} m_K (1 + \beta \Delta t) v_K^{n+1} = \sum_{K \in \mathcal{V}} m_K v_K^n + \alpha \Delta t \sum_{K \in \mathcal{V}} m_K u_K^n, \quad \forall n \in \{0, \dots, N\}. \quad (4.47)$$

Assume that $v_{K_\star^n}^n > 0$ or $u_{K_\star^n}^n > 0$ for some $K_\star^n \in \mathcal{V}$, as this is the case for $n = 0$ because of the assumption on the initial data u_0 and v_0 , then we deduce from (4.47) and from the non-negativity of v_K^n and u_K^n proved in Lemma 4.6 that

$$\sum_{K \in \mathcal{V}} m_K (1 + \beta \Delta t) v_K^{n+1} > 0.$$

In particular, there exists $K_\star^{n+1} \in \mathcal{V}$ such that $v_{K_\star^{n+1}}^{n+1}$ is (strictly) positive and

$$\sum_{K \in \mathcal{V}} m_K v_K^{n+1} := \kappa_{n+1} > 0.$$

One concludes the proof by setting $\kappa = \min_{n=1, \dots, N+1} \kappa_n$. \square

We give now the definition of D -transmissive path, which was introduced in [16, Definition 3.4].

Definition 4.11. A D -transmissive path w joining $K_i \in \mathcal{V}$ to $K_f \in \mathcal{V}$ consists in a list of vertices $(K_q)_{0 \leq q \leq M}$ such that $K_i = K_0$, $K_f = K_M$, with $K_q \neq K_\ell$ if $q \neq \ell$, and such that $\sigma_{K_q K_{q+1}} \in \mathcal{E}$ with $\bar{D}_{K_q K_{q+1}} > 0$ for all $q \in \{0, \dots, M-1\}$. We denote by $\mathcal{W}(K_i, K_f)$ the set of the transmissive path joining $K_i \in \mathcal{V}$ to $K_f \in \mathcal{V}$.

We now state a result which is proved in [16, Lemma 3.5].

Lemma 4.12. For all $(K_i, K_f) \in \mathcal{V}^2$ there exists a transmissive path $w \in \mathcal{W}(K_i, K_f)$.

We have now introduced all the necessary tools for proving the main result of this section.

Lemma 4.13. Assume that $\int_{\Omega} u_0(\mathbf{x}) d\mathbf{x} > 0$ or $\int_{\Omega} v_0(\mathbf{x}) d\mathbf{x} > 0$, then there exists $r_h > 0$ depending on the data as well as on the mesh \mathcal{T} and Δt such that

$$v_K^{n+1} \geq r_h, \quad \forall K \in \mathcal{V}, \forall n \in \{0, \dots, N\}. \quad (4.48)$$

Proof. Thanks to Lemma 4.10, we know that there exists K_i such that $v_{K_i}^{n+1} > 0$. Let $K_f \in \mathcal{V}$, then there exists a D -transmissive path $w = (K_q)_{0 \leq q \leq M} \in \mathcal{W}(K_i, K_f)$ thanks to Lemma 4.12, with $K_0 = K_i$ and $K_M = K_f$.

Thanks to Lemmas 4.5 and 4.7, we know that there exists C such that

$$\sum_{n=0}^N \Delta t \sum_{\sigma_{KL} \in \mathcal{E}} |D_{KL}| \eta_{KL}^{n+1} (p(v_K^{n+1}) - p(v_L^{n+1}))^2 \leq C.$$

In particular, this ensures that

$$D_{K_q K_{q+1}} \eta_{K_q K_{q+1}}^{n+1} \left(p(v_{K_q}^{n+1}) - p(v_{K_{q+1}}^{n+1}) \right)^2 \leq \frac{C}{\Delta t}, \quad \forall q \in \{0, \dots, M-1\}.$$

Assume now that $v_{K_q}^{n+1} > 0$, as this is the case for $q = 0$, then $\eta_{K_q K_{q+1}}^{n+1} \geq \eta_{K_q}^{n+1} > 0$. Then one has

$$\left(p(v_{K_q}^{n+1}) - p(v_{K_{q+1}}^{n+1}) \right)^2 \leq \frac{C}{\Delta t D_{K_q K_{q+1}} \eta_{K_q K_{q+1}}^{n+1}} < \infty. \quad (4.49)$$

Since $\lim_{v \rightarrow 0} p(v) = -\infty$, we deduce from (4.49) that $p(v_{K_{q+1}}^{n+1}) > -\infty$, hence $v_{K_{q+1}}^{n+1} > 0$. A straightforward induction provides that $v_{K_f}^{n+1} > 0$, and since K_f was chosen arbitrarily, we obtain that

$$v_K^{n+1} > 0, \quad \forall K \in \mathcal{V}.$$

Since the set $\mathcal{V} \times \{0, \dots, N\}$ is finite, we can conclude that there exists r_h such that (4.48) holds. \square

4.5.5 Existence of a discrete solution

Proposition 4.14. Given $(u_K^n, v_K^n)_{K \in \mathcal{V}}$ such that $u_{\mathcal{M}, \Delta t}(\cdot, n\Delta t)$ and $v_{\mathcal{M}, \Delta t}(\cdot, n\Delta t)$ are nonnegative, then there exists (at least) one solution $(u_K^{n+1}, v_K^{n+1})_{K \in \mathcal{V}}$ of the scheme (4.13), (4.21). Moreover, $u_{\mathcal{M}, \Delta t}(\cdot, (n+1)\Delta t)$ and $v_{\mathcal{M}, \Delta t}(\cdot, (n+1)\Delta t)$ are nonnegative.

Proof. The case where $(u_K^n, v_K^n)_{K \in \mathcal{V}} \equiv 0$ has to be treated apart. In this very particular case, it is easy to check that $(u_K^{n+1}, v_K^{n+1})_{K \in \mathcal{V}} \equiv 0$ is a solution to the scheme.

Let us now focus on the case where u_K^n or v_K^n is strictly positive for some $K \in \mathcal{V}$. Because of the weak coupling on the numerical scheme, we can first solve (4.21), and afterwards (4.13). The existence of a solution $(v_K^{n+1})_{K \in \mathcal{V}}$ can be proved by slightly adapting the proof of [16, Proposition 3.11], which relies on a topological argument [54, 23]. The main difficulty comes from the fact the scheme (4.21) is not continuous w.r.t. $(v_K^{n+1})_{K \in \mathcal{V}}$ on $(\mathbb{R}_+)^{\#\mathcal{V}}$, but Lemma 4.13 ensures that no component v_K^{n+1} of the discrete solution can go close to 0. Let us detail now the proof.

Let $\mu \in [0, 1]$, we denote by $(v_{K,\mu}^{n+1})_{K \in \mathcal{V}}$ the solution (if it exists) to the numerical scheme

$$\begin{aligned} \frac{v_{K,\mu}^{n+1} - v_K^n}{\Delta t} m_K + \mu \sum_{\sigma_{KL} \in \mathcal{E}_K} D_{KL} \eta_{KL,\mu}^{n+1} (p(v_{K,\mu}^{n+1}) - p(v_{L,\mu}^{n+1})) \\ + (1 - \mu) \sum_{\sigma_{KL} \in \mathcal{E}_K} |D_{KL}| (p(v_{K,\mu}^{n+1}) - p(v_{L,\mu}^{n+1})) = \alpha u_K^n m_K - \beta v_{K,\mu}^{n+1} m_K. \end{aligned} \quad (4.50)$$

In the above scheme, we have set

$$\eta_{KL,\mu}^{n+1} = \begin{cases} \max_{v \in J_{KL,\mu}^{n+1}} \eta(v) & \text{if } D_{KL} \geq 0, \\ \min_{v \in J_{KL,\mu}^{n+1}} \eta(v) & \text{if } D_{KL} < 0, \end{cases}$$

where $J_{KL,\mu}^{n+1} = [\min(v_{K,\mu}^{n+1}, v_{L,\mu}^{n+1}), \max(v_{K,\mu}^{n+1}, v_{L,\mu}^{n+1})]$. Reproducing carefully the analysis carried out in §4.5.2 and §4.5.4, we get that for all $\mu \in [0, 1]$,

$$\sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \left(\phi(v_{K,\mu}^{n+1}) - \phi(v_{L,\mu}^{n+1}) \right)^2 \leq \sum_{\sigma_{KL} \in \mathcal{E}} D_{KL} \eta_{KL,\mu}^{n+1} \left(p(v_{K,\mu}^{n+1}) - p(v_{L,\mu}^{n+1}) \right)^2 \leq C \quad (4.51)$$

and, that there exists $\epsilon > 0$ such that

$$v_{K,\mu}^{n+1} \geq \epsilon > 0, \quad \forall K \in \mathcal{V}. \quad (4.52)$$

This ensures in particular that for all $\mu \in [0, 1]$, the solutions of (4.50) stay in the interior of a compact subset \mathcal{K} of $\mathbb{R}^{\#\mathcal{V}}$ such that

$$\text{dist} \left(\mathcal{K}, (\mathbb{R}_-)^{\#\mathcal{V}} \right) \geq \frac{\epsilon}{2}.$$

Define the function $\Upsilon : \mathcal{K} \times [0, 1] \rightarrow \mathbb{R}^{\#\mathcal{V}}$ by : $\forall K \in \mathcal{V}$,

$$\begin{aligned} \Upsilon_K((w_K)_K, \mu) = \frac{w_K - v_K^n}{\Delta t} m_K + \mu \sum_{\sigma_{KL} \in \mathcal{E}_K} D_{KL} \eta_{KL,\mu}^{n+1} (p(w_K) - p(w_L)) \\ + (1 - \mu) \sum_{\sigma_{KL} \in \mathcal{E}_K} |D_{KL}| (p(w_K) - p(w_L)) - \alpha u_K^n m_K + \beta w_K m_K. \end{aligned}$$

The function Υ is uniformly continuous on $\mathcal{K} \times [0, 1]$, and for all $\mu \in [0, 1]$ the solution $v_{K,\mu}^{n+1}$ of the nonlinear system

$$\Upsilon \left((v_{K,\mu}^{n+1})_{K \in \mathcal{V}}, \mu \right) = 0 \quad (4.53)$$

cannot reach $\partial \mathcal{K}$. For $\mu = 0$, the system is monotone, so that the system (4.53) admits a unique solution, whose topological degree is equal to 1 (we refer to [30, Proposition 3.1] for a proof of this property). The topological degree being constant w.r.t. $\mu \in [0, 1]$, the system (4.53) admits at least one solution for $\mu = 1$, concluding the proof of the existence of $(v_K^{n+1})_{K \in \mathcal{V}}$.

The existence proof for $(u_K^{n+1})_{K \in \mathcal{V}}$ is similar but simpler since

1. the *a priori* estimate $0 \leq u_K^{n+1} \leq 1$ is sufficient for the claim, and no energy estimate is needed here ;
2. the scheme (4.13) depends in a uniformly continuous way on $(u_K^{n+1})_{K \in \mathcal{V}}$ on the compact subset $[-1, 2]^{\#\mathcal{V}}$ of $\mathbb{R}^{\#\mathcal{V}}$.

Therefore, we let to the reader the care of checking the proof for self-conviction. \square

4.6 Compactness estimates on the family of discrete solutions.

In this section, we derive estimates on differences of time and space translates of the discrete solutions necessary to prove the relative compactness property of the sequence of approximate solutions. To do this, we use the increasing function $\xi : [0, 1] \rightarrow \mathbb{R}$ and $\phi : [0, +\infty) \rightarrow \mathbb{R}$ defined in equation (4.5).

For all $K \in \mathcal{V}_h$ and for all $n \geq 1$, we denote by $\xi_K^n = \xi(u_K^n)$, and by $\xi_{\mathcal{T}_h, \Delta t_h}$ the corresponding piecewise affine in space and constant in time reconstruction in $\mathcal{H}_{\mathcal{T}_h, \Delta t_h}$, and by $\xi_{\mathcal{M}_h, \Delta t_h}$ the piecewise constant reconstruction in $\mathcal{X}_{\mathcal{M}_h, \Delta t_h}$.

4.6.1 Time translate estimate.

We give below the time translate estimate for the family $(\xi_{\mathcal{M}_h, \Delta t_h})_h$ expressed using the function ξ defined in (4.5). We denote by $Q_{t_f - \tau} = \Omega \times (0, t_f - \tau)$, for all $\tau \in (0, t_f)$.

Lemma 4.15. *There exists a constant $C_{\xi, t}$ and $C_{v, t}$ independent of h and τ such that,*

$$\iint_{Q_{t_f - \tau}} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - \xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt \leq C_{\xi, t} (\tau + \Delta t_h), \quad (4.54)$$

$$\iint_{Q_{t_f - \tau}} |v_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - v_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt \leq C_{v, t} (\tau + \Delta t_h), \quad (4.55)$$

for all $\tau \in (0, t_f)$.

Proof. We consider the quantity $A_h(t)$ defined by

$$A_h(t) = \int_{\Omega} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - \xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x}, \quad \text{for all } t \in (0, t_f - \tau),$$

which implies, that

$$\iint_{Q_{t_f - \tau}} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - \xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt = \int_0^{t_f - \tau} A_h(t) dt.$$

For $t \in (0, t_f]$, we denote by $\nu(t) \in \{0, \dots, N\}$ the unique positive integer such that $t_{\nu(t)} < t \leq t_{\nu(t)+1}$, so that, we can rewrite $A_h(t)$ as

$$A_h(t) = \sum_{K \in \mathcal{V}_h} \left(\xi_K^{\nu(t+\tau)+1} - \xi_K^{\nu(t)+1} \right)^2 m_K, \quad \text{for all } t \in (0, t_f - \tau).$$

We observe, using the definition (4.5) of the function ξ , that

$$\begin{aligned} \left(\xi_K^{\nu(t+\tau)+1} - \xi_K^{\nu(t)+1} \right)^2 &\leq C \left(u_K^{\nu(t+\tau)+1} - u_K^{\nu(t)+1} \right) \times \left(\xi(u_K^{\nu(t+\tau)+1}) - \xi(u_K^{\nu(t)+1}) \right) \\ &= C \sum_{n=\nu(t)+1}^{\nu(t+\tau)} (u_K^{n+1} - u_K^n) \left(\xi(u_K^{\nu(t+\tau)+1}) - \xi(u_K^{\nu(t)+1}) \right). \end{aligned}$$

Now, using scheme (4.13) then gathering by edges and using the weighted Young inequality, we obtain

$$\begin{aligned}
A_h(t) &\leq C \sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} \left[\left(\Lambda_{KL} a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) - \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}) \right) \right. \\
&\quad \times \left(\left(\xi(u_K^{\nu(t)+1}) - \xi(u_L^{\nu(t)+1}) \right) - \left(\xi(u_K^{\nu(t+\tau)+1}) - \xi(u_L^{\nu(t+\tau)+1}) \right) \right) \Big] \\
&\quad + C \sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h \sum_{K \in \mathcal{V}_h} m_K f(u_K^{n+1}) \left(\xi(u_K^{\nu(t+\tau)+1}) - \xi(u_K^{\nu(t)+1}) \right) \\
&\leq C (A_{1,h}(t) + A_{2,h}(t) + A_{3,h}(t) + A_{4,h}(t) + A_{5,h}(t)),
\end{aligned}$$

where, we have set

$$\begin{aligned}
A_{1,h}(t) &= \sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} |\Lambda_{KL}| a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2, \\
A_{2,h}(t) &= \frac{\|a\|_\infty^2 + \|\chi\|_\infty^2}{2} \sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} |\Lambda_{KL}| \left(\xi(u_K^{\nu(t)+1}) - \xi(u_L^{\nu(t)+1}) \right)^2, \\
A_{3,h}(t) &= \frac{\|a\|_\infty^2 + \|\chi\|_\infty^2}{2} \sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} |\Lambda_{KL}| \left(\xi(u_K^{\nu(t+\tau)+1}) - \xi(u_L^{\nu(t+\tau)+1}) \right)^2, \\
A_{4,h}(t) &= \sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} |\Lambda_{KL}| (v_K^{n+1} - v_L^{n+1})^2, \\
A_{5,h}(t) &= \sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h \sum_{K \in \mathcal{V}_h} m_K |f(u_K^{n+1})| \left| \xi(u_K^{\nu(t+\tau)+1}) - \xi(u_K^{\nu(t)+1}) \right|.
\end{aligned}$$

Now, we introduce the characteristic function $\rho(n, t, \tau) = 1$ if $t < n\Delta t_h \leq t + \tau$ and $\rho(n, t, \tau) = 0$ otherwise. Let $(a^n)_{n \in \{0, \dots, N\}}$ be a family of non negative real values, we have the following properties

$$\int_0^{t_f - \tau} \rho(n, t, \tau) dt = \int_{t_n - \tau}^{t_n} dt = \tau, \quad \sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h = \sum_{n; t < t_n \leq t + \tau} t_{n+1} - t_n \leq \tau + \Delta t_h,$$

and

$$\int_0^{t_f - \tau} \sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h a^{n+1} dt = \int_0^{t_f - \tau} \sum_{n=0}^N \Delta t_h a^{n+1} \rho(n, t, \tau) dt = \tau \sum_{n=0}^N \Delta t_h a^{n+1}.$$

Using these properties, Lemma 4.5, the a priori estimates (4.40)–(4.41) on u_K^{n+1} and v_K^{n+1} , and the L^∞ bound of the function ξ , one can deduce that

$$\int_0^{t_f - \tau} A_{1,h}(t) dt \leq C(\tau + \Delta t_h), \quad \int_0^{t_f - \tau} A_{4,h}(t) dt \leq C(\tau + \Delta t_h),$$

and

$$\int_0^{t_f - \tau} A_{5,h}(t) dt \leq C(\tau + \Delta t_h),$$

for some constant $C > 0$.

It remains to show that $\int_0^{t_f-\tau} A_{2,h}(t) dt \leq C(\tau + \Delta t_h)$, and $\int_0^{t_f-\tau} A_{3,h}(t) dt \leq C(\tau + \Delta t_h)$, for some constant $C > 0$.

Consider the function ζ defined by $\zeta(n, t) = 1$ if $\nu(t) = n$ and $\zeta(n, t) = 0$ otherwise. We have the following result (see ...)

$$\begin{aligned} \int_0^{t_f-\tau} \left(\sum_{n=\nu(t)+1}^{\nu(t+\tau)} \Delta t_h \right) a^{\nu(t+\tau)+1} dt &\leq (\tau + \Delta t_h) \int_0^{t_f-\tau} \sum_{m=0}^N a^{m+1} \zeta(m, t + \tau) dt \\ &= (\tau + \Delta t_h) \sum_{m=0}^N a^{m+1} \int_{t^m-\tau}^{t^{m+1}-\tau} dt = (\tau + \Delta t_h) \sum_{m=0}^N a^{m+1} \Delta t_h. \end{aligned}$$

One can conclude the proof using this property, Lemma 4.3, and Lemma 4.5. The proof of estimate (4.55) is similar. \square

We now extend by zero the functions $\xi_{\mathcal{M}_h, \Delta t_h}$ and $\xi_{\mathcal{T}_h, \Delta t_h}$ outside of Q_{t_f} and give the time translate estimate over \mathbb{R}^3 for $\xi_{\mathcal{M}_h, \Delta t_h}$. Indeed, there exists a constant $C > 0$ independent of h and τ such that

$$\int_{\mathbb{R}} \int_{\mathbb{R}^2} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - \xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt \leq C(\tau + \Delta t_h), \quad \text{for all } \tau \in (0, t_f).$$

Proof. Using the extension by zero of $\xi_{\mathcal{M}_h, \Delta t_h}$ outside of Q_{t_f} , one has

$$\begin{aligned} &\int_{\mathbb{R}} \int_{\mathbb{R}^2} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - \xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt \\ &= \iint_{Q_{t_f-\tau}} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t + \tau) - \xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt + \int_{t_f-\tau}^{t_f} \int_{\Omega} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)|^2 d\mathbf{x} dt. \end{aligned}$$

One can deduce the proof using Lemma 4.15 and the L^∞ bound of the function ξ . \square

We give now the space translate estimate on $\xi_{\mathcal{M}_h, \Delta t_h}$.

4.6.2 Space translate estimate.

Lemma 4.16. *There exists a constant $C_{\xi,s}$ independent of h and \mathbf{y} such that,*

$$\int_0^{t_f} \int_{\mathbb{R}^2} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x} + \mathbf{y}, t) - \xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)| d\mathbf{x} dt \leq C_{\xi,s} (|\mathbf{y}| + h), \quad (4.56)$$

$$\int_0^{t_f} \int_{\mathbb{R}^2} |v_{\mathcal{M}_h, \Delta t_h}(\mathbf{x} + \mathbf{y}, t) - v_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)| d\mathbf{x} dt \leq C_{\xi,s} (|\mathbf{y}| + h), \quad (4.57)$$

for all $\mathbf{y} \in \mathbb{R}^2$.

Proof. We follow the same proof used in chapter 3. We first prove that

$$\int_0^{t_f} \int_{\mathbb{R}^2} |\xi_{\mathcal{T}_h, \Delta t_h}(\mathbf{x} + \mathbf{y}, t) - \xi_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t)| d\mathbf{x} dt \leq C|\mathbf{y}|, \quad \text{for all } \mathbf{y} \in \mathbb{R}^2,$$

then, using the extension by zero of $\xi_{\mathcal{M}_h, \Delta t_h}$ outside of Q_{t_f} and the triangle inequality, we get

$$\int_0^{t_f} \int_{\mathbb{R}^2} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x} + \mathbf{y}, t) - \xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)| d\mathbf{x} dt \leq A + B,$$

where

$$A = \int_0^{t_f} \int_{\mathbb{R}^2} |\xi_{\mathcal{T}_h, \Delta t_h}(\mathbf{x} + \mathbf{y}, t) - \xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t)| d\mathbf{x} dt,$$

$$B = 2 \iint_{Q_{t_f}} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t) - \xi_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t)| d\mathbf{x} dt.$$

One can conclude the proof using the following estimate

$$\iint_{Q_{t_f}} |\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t) - \xi_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t)| d\mathbf{x} dt \leq Ch. \quad (4.58)$$

The proof of estimate (4.57) follows the same lines. This ends the proof of the lemma. \square

4.7 Convergence

Lemma 4.17 (Strong convergence in $L^1(Q_{t_f})$). *There exists a subsequence of the sequence $(\xi(u)_{\mathcal{M}_h, \Delta t_h})_{h>0}$ (resp. $(v_{\mathcal{M}_h, \Delta t_h})_{h>0}$) which converges in $L^1(Q_{t_f})$ to some function $\xi(u) \in L^2(0, t_f; H^1(\Omega))$ (resp. $v \in L^2(0, t_f; H^1(\Omega))$).*

Proof. Let us consider the function $\xi_{\mathcal{M}_h, \Delta t_h}$ on Q_{t_f} and its extension by zero outside of Q_{t_f} . Lemma 4.15, Lemma 4.16, and the boundedness of ξ due to Lemma 4.5.1 ensure that the sequence $(\xi(u)_{\mathcal{M}_h, \Delta t_h})_{h>0}$ verifies the assumptions of the Kolmogorov compactness criterion (see e.g. [38, 32]), where the third item of the theorem is satisfied using the triangle inequality : for any $\eta \in \mathbb{R}^2$ and $\tau \in \mathbb{R}$,

$$\begin{aligned} \|\xi_{\mathcal{M}_h, \Delta t_h}(\cdot + \eta, \cdot + \tau) - \xi_{\mathcal{M}_h, \Delta t_h}(\cdot, \cdot)\|_{L^1(\mathbb{R}^3)} &\leq \|\xi_{\mathcal{M}_h, \Delta t_h}(\cdot + \eta, \cdot) - \xi_{\mathcal{M}_h, \Delta t_h}(\cdot, \cdot)\|_{L^1(\mathbb{R}^3)} \\ &\quad + \|\xi_{\mathcal{M}_h, \Delta t_h}(\cdot, \cdot + \tau) - \xi_{\mathcal{M}_h, \Delta t_h}(\cdot, \cdot)\|_{L^1(\mathbb{R}^3)}. \end{aligned}$$

Kolmogorov's theorem ensures that the sequence $(\xi(u)_{\mathcal{M}_h, \Delta t_h})_{h>0}$ is relatively compact in $L^1(Q_{t_f})$, that implies the existence of a subsequence of $(\xi(u)_{\mathcal{M}_h, \Delta t_h})_{h>0}$ such that

$$\xi(u)_{\mathcal{M}_h, \Delta t_h} \longrightarrow \xi^* \quad \text{strongly in } L^1(Q_{t_f}). \quad (4.59)$$

Furthermore, as ξ is a continuous and nondecreasing function on $[0, 1]$, there exists a unique $u(\mathbf{x}, t)$ defined by

$$u(\mathbf{x}, t) = \xi^{-1}(\xi^*(\mathbf{x}, t)), \quad \text{for a.e. } (\mathbf{x}, t) \in Q_{t_f - \tau}.$$

Since ξ^{-1} is well defined and continuous, applying the L^∞ bound on $u_{\mathcal{M}_h, \Delta t_h}$ and the dominated convergence theorem to $u_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t) = \xi^{-1}(\xi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t))$, $\forall (\mathbf{x}, t) \in \omega_K \times (0, t_f)$, $\forall K \in \mathcal{V}$, we get

$$u_{\mathcal{M}_h, \Delta t_h} \longrightarrow u \text{ a.e. in } Q_{t_f} \text{ and strongly in } L^p(Q_{t_f}) \text{ for } p < +\infty.$$

It remains to show that $\xi(u) \in L^2(0, t_f; H^1(\Omega))$. Indeed, plugging estimate (4.41) into estimate (4.22), one gets that $\nabla \xi(u)_{\mathcal{T}_h, \Delta t_h}$ is uniformly bounded in $(L^2(Q_{t_f}))^2$. It follows that the sequence $(\xi(u)_{\mathcal{T}_h, \Delta t_h})_h$ is bounded in $L^2(0, t_f; H^1(\Omega))$ since $\xi(u)_{\mathcal{T}_h, \Delta t_h}$ is uniformly bounded in $L^2(Q_{t_f})$. Consequently, the sequence $(\xi(u)_{\mathcal{T}_h, \Delta t_h})_h$ converges weakly, up to an unlabeled subsequence, to a function $\tilde{\xi}$ in $L^2(0, t_f; H^1(\Omega))$.

According to estimate (4.58), the sequences $(\xi(u)_{\mathcal{T}_h, \Delta t_h})_h$ and $(\xi(u)_{\mathcal{M}_h, \Delta t_h})_h$ have the same limit, as a consequence $\xi(u) = \xi^* = \tilde{\xi} \in L^2(0, t_f; H^1(Q_{t_f}))$. \square

4.7.1 Identification as a weak solution

It remains to be shown that (u, v) is a weak solution to the continuous problem (4.1)–(4.3) in the sense of Definition 4.1. To do this, we consider a test function $\psi \in \mathcal{D}(\bar{\Omega} \times [0, t_f])$, and denote by $\psi_K^n = \psi(\mathbf{x}_K, t_n)$, for all $K \in \mathcal{V}_h$ and all $n \in \{0, \dots, N_h\}$. Let us focus on the convergence of the first equation of scheme (4.13)–(4.21), the convergence of the second equation being similar.

Multiplying the first equation (4.13) by $\Delta t_h \psi_K^n$ and summing over $n \in \{0, \dots, N_h\}$ and $K \in \mathcal{V}_h$ yields, after a reorganization of the sum (see e.g. [34]),

$$\mathcal{A}_h + \mathcal{B}_h + \mathcal{C}_h + \mathcal{D}_h = \mathcal{F}_h, \quad (4.60)$$

where

$$\begin{aligned} \mathcal{A}_h &= \sum_{n=0}^{N_h} \sum_{K \in \mathcal{V}_h} (u_K^{n+1} - u_K^n) \psi_K^n m_K, & \mathcal{F}_h &= \sum_{n=0}^{N_h} \Delta t_h \sum_{K \in \mathcal{V}_h} f(u_K^{n+1}) \psi_K^n m_K, \\ \mathcal{B}_h &= \sum_{n=0}^{N_h} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} \Lambda_{KL} \left(a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1}) - \sqrt{a_{KL}^{n+1}} (\xi(u_K^{n+1}) - \xi(u_L^{n+1})) \right) (\psi_K^n - \psi_L^n), \\ \mathcal{C}_h &= \sum_{n=0}^{N_h} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} \Lambda_{KL} \sqrt{a_{KL}^{n+1}} (\xi(u_K^{n+1}) - \xi(u_L^{n+1})) (\psi_K^n - \psi_L^n), \\ \mathcal{D}_h &= - \sum_{n=0}^{N_h} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} \Lambda_{KL} \mu_{KL}^{n+1} a_{KL}^{n+1} (v_K^{n+1} - v_L^{n+1}) (\psi_K^n - \psi_L^n). \end{aligned}$$

Time evolution term Note that $\psi_K^{N_h+1} = 0$ for all $K \in \mathcal{V}_h$, then, performing summation by parts in time, the term \mathcal{A}_h can be rewritten

$$\begin{aligned} \mathcal{A}_h &= \sum_{n=0}^{N_h} \sum_{K \in \mathcal{V}_h} u_K^{n+1} \psi_K^n m_K - \sum_{n=1}^{N_h} \sum_{K \in \mathcal{V}_h} u_K^n \psi_K^n m_K - \sum_{K \in \mathcal{V}_h} u_K^0 \psi_K^0 m_K \\ &= - \sum_{n=0}^{N_h} \Delta t_h \sum_{K \in \mathcal{V}_h} u_K^{n+1} \frac{\psi_K^{n+1} - \psi_K^n}{\Delta t_h} m_K - \sum_{K \in \mathcal{V}_h} u_K^0 \psi_K^0 m_K \\ &= - \iint_{Q_{t_f}} u_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t) \partial_t \psi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, t) \, d\mathbf{x} \, dt - \int_{\Omega} u_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, 0) \psi_{\mathcal{M}_h, \Delta t_h}(\mathbf{x}, 0) \, d\mathbf{x}. \end{aligned}$$

Thanks to the regularity of ψ , and the uniform convergence in $L^q(Q_{t_f})$, for all $q \in [1, \infty)$, of the sequence $(u_{\mathcal{M}_h, \Delta t_h})_h$ towards u , it follows that

$$\mathcal{A}_h \longrightarrow - \iint_{Q_{t_f}} u(\mathbf{x}, t) \partial_t \psi(\mathbf{x}, t) \, d\mathbf{x} \, dt - \int_{\Omega} u(\mathbf{x}, 0) \psi(\mathbf{x}, 0) \, d\mathbf{x}, \quad \text{as } h \rightarrow 0.$$

Diffusion term Let us first prove that $\lim_{h \rightarrow 0} \mathcal{B}_h = 0$.

For all $\sigma_{KL} \in \mathcal{E}_h$ and all $n \in \{0, \dots, N_h\}$, we denote by \bar{a}_{KL}^{n+1} the quantity defined by

$$\bar{a}_{KL}^{n+1} = \begin{cases} \left(\frac{\xi(u_K^{n+1}) - \xi(u_L^{n+1})}{u_K^{n+1} - u_L^{n+1}} \right)^2 & \text{if } u_K^{n+1} \neq u_L^{n+1}, \\ a(u_K^{n+1}) & \text{if } u_K^{n+1} = u_L^{n+1}. \end{cases}$$

Then, the term \mathcal{B}_h rewrites

$$\mathcal{B}_h = \sum_{n=0}^{N_h} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} \Lambda_{KL} \sqrt{a_{KL}^{n+1}} \left(\sqrt{a_{KL}^{n+1}} - \sqrt{\bar{a}_{KL}^{n+1}} \right) (u_K^{n+1} - u_L^{n+1}) (\psi_K^n - \psi_L^n).$$

Now, using the Cauchy-Schwarz inequality, we get

$$|\mathcal{B}_h| \leq \left(\sum_{n=0}^{N_h} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} |\Lambda_{KL}| a_{KL}^{n+1} (u_K^{n+1} - u_L^{n+1})^2 \right)^{\frac{1}{2}} \times \mathcal{R}_h^{\frac{1}{2}},$$

where, \mathcal{R}_h is given by

$$\mathcal{R}_h = \sum_{n=0}^{N_h} \Delta t_h \sum_{\sigma_{KL} \in \mathcal{E}_h} |\Lambda_{KL}| \left(\sqrt{a_{KL}^{n+1}} - \sqrt{\bar{a}_{KL}^{n+1}} \right)^2 (\psi_K^n - \psi_L^n)^2.$$

Using Lemma 4.5 and Proposition 4.9, one has $|\mathcal{B}_h| \leq C \mathcal{R}_h^{\frac{1}{2}}$. However, in order to prove that $\lim_{h \rightarrow 0} \mathcal{B}_h = 0$, it suffices to prove that $\lim_{h \rightarrow 0} \mathcal{R}_h = 0$.

For all $T \in \mathcal{T}_h$, we denote by

$$\bar{\xi}_T^{n+1} = \max_{\mathbf{x} \in T} \left(\xi(p)_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t_{n+1}) \right), \quad \underline{\xi}_T^{n+1} = \min_{\mathbf{x} \in T} \left(\xi(p)_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t_{n+1}) \right),$$

and for all $(\mathbf{x}, t) \in T \times (t_n, t_{n+1})$, by

$$\bar{\xi}_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) = \bar{\xi}_T^{n+1}, \quad \underline{\xi}_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) = \underline{\xi}_T^{n+1}.$$

Consider the uniform continuous function $\sqrt{a \circ \xi^{-1}}$ defined on the closed and bounded interval $[0, \xi(1)]$, and let η be its modulus of continuity, then we have

$$\left| \sqrt{a_{KL}^{n+1}} - \sqrt{\bar{a}_{KL}^{n+1}} \right| \leq \eta \left(\bar{\xi}_T^{n+1} - \underline{\xi}_T^{n+1} \right), \quad \text{for all } \sigma_{KL} \in \mathcal{E}_T. \quad (4.61)$$

Therefore, using this inequality in the definition of \mathcal{R}_h , we get

$$0 \leq \mathcal{R}_h \leq \mathcal{Q}_h \quad (4.62)$$

where,

$$\mathcal{Q}_h = \sum_{n=0}^{N_h} \Delta t_h \sum_{T \in \mathcal{T}_h} \left(\eta \left(\bar{\xi}_T^{n+1} - \underline{\xi}_T^{n+1} \right) \right)^2 \sum_{\sigma_{KL} \in \mathcal{E}_T} |\lambda_{KL}^T| (\psi_K^n - \psi_L^n)^2. \quad (4.63)$$

Thanks to Lemma 4.4, one can deduce that the inequality (4.62) implies that

$$0 \leq \mathcal{R}_h \leq C \iint_{Q_{t_f}} \eta \left(\bar{\xi}_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) - \underline{\xi}_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) \right) d\mathbf{x} dt,$$

where C is independent of h , and Δt_h .

Therefore, it suffices to show that $\bar{\xi}_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) - \underline{\xi}_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) \rightarrow 0$ for a.e. in Q_{t_f} to consequently prove that $\lim_{h \rightarrow 0} \mathcal{R}_h = 0$. However, by a simple generalization of Lemma A.1 and by the help of Lemma 4.3 and Proposition 4.9, it follows that

$$\iint_{Q_{t_f}} \left| \bar{\xi}_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) - \underline{\xi}_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) \right| d\mathbf{x} dt \leq Ch \left(\iint_{Q_{t_f}} \left| \nabla \xi(u)_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) \right|^2 d\mathbf{x} dt \right)^{\frac{1}{2}} \leq Ch.$$

As a consequence, up to a subsequence, one has

$$\lim_{h \rightarrow 0} \mathcal{B}_h = \lim_{h \rightarrow 0} \mathcal{R}_h = \lim_{h \rightarrow 0} \mathcal{Q}_h = 0.$$

We now focus on the term \mathcal{C}_h and prove that

$$\lim_{h \rightarrow 0} \mathcal{C}_h = \iint_{Q_{t_f}} \Lambda(\mathbf{x}) a(u) \nabla u \cdot \nabla \psi d\mathbf{x} dt.$$

To do this, we introduce the term \mathcal{C}'_h defined by

$$\mathcal{C}'_h := \iint_{Q_{t_f}} \Theta_{\mathcal{T}_h, \Delta t_h} \Lambda(\mathbf{x}) \nabla \xi(u)_{\mathcal{T}_h, \Delta t_h} \cdot \nabla \psi_{\mathcal{T}_h, \Delta t_h}(\cdot, t - \Delta t_h) d\mathbf{x} dt,$$

where $\Theta_{\mathcal{T}_h, \Delta t_h}$ is a continuous function given by

$$\Theta_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) = \sqrt{a \circ \xi^{-1}}(\Upsilon_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t)), \quad \forall \mathbf{x} \in T, \forall t \in (t_n, t_{n+1}], \forall T \in \mathcal{T}_h,$$

and $\Upsilon_{\mathcal{T}_h, \Delta t_h}$ is a function defined by

$$\Upsilon_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) = \xi(u)_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}_T, t), \quad \forall \mathbf{x} \in T, \forall t \in (t_n, t_{n+1}], \forall T \in \mathcal{T}_h.$$

Using again a generalization of Lemma A.1 as well as the boundedness of the continuous function $\sqrt{a \circ \xi^{-1}}$, we obtain

$$\begin{aligned} \Upsilon_{\mathcal{T}_h, \Delta t_h} &\longrightarrow \xi(u) && \text{in } L^2(Q_{t_f}) \text{ as } h \rightarrow 0, \\ \Theta_{\mathcal{T}_h, \Delta t_h} &\longrightarrow \sqrt{a(u)} && \text{in } L^2(Q_{t_f}) \text{ as } h \rightarrow 0. \end{aligned} \quad (4.64)$$

It remains to verify that $|\mathcal{C}_h - \mathcal{C}'_h| \longrightarrow 0$, when h tends to zero.

We denote by

$$a_T^{n+1} = (\Theta_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}_T, t_{n+1}))^2, \quad \forall T \in \mathcal{T}_h, \forall n \in \{0, \dots, N_h\}.$$

The discretization of the term \mathcal{D}'_m is written as

$$\mathcal{C}'_m = \sum_{n=0}^{N_h} \Delta t_h \sum_{T \in \mathcal{T}_h} \sqrt{a_T^{n+1}} \sum_{\sigma_{KL} \in \mathcal{E}_T} \lambda_{KL}^T (\xi(u_K^{n+1}) - \xi(u_L^{n+1})) (\psi_K^n - \psi_L^n).$$

Performing the same way as for the inequality (4.61), one has

$$\left| \sqrt{a_{KL}^{n+1}} - \sqrt{a_T^{n+1}} \right| \leq \eta \left(\bar{\xi}_T^{n+1} - \underline{\xi}_T^{n+1} \right), \quad \text{for all } \sigma_{KL} \in \mathcal{E}_T.$$

Therefore, using the Cauchy-Schwarz inequality, Lemma 4.3, Lemma 4.4, and Proposition 4.9, we deduce that there exists a constant C does not depend on h such that

$$\begin{aligned} |\mathcal{C}_m - \mathcal{C}'_m|^2 &\leq \left(\sum_{n=0}^{N_h} \Delta t_h \sum_{T \in \mathcal{T}_h} \eta \left(\bar{\xi}_T^{n+1} - \underline{\xi}_T^{n+1} \right) \sum_{\sigma_{KL} \in \mathcal{E}_T} |\lambda_{KL}^T| |\xi(u_K^{n+1}) - \xi(u_L^{n+1})| |\psi_K^n - \psi_L^n| \right)^2 \\ &\leq \mathcal{Q}_h \times \sum_{n=0}^{N_h} \Delta t_h \sum_{T \in \mathcal{T}_h} \sum_{\sigma_{KL} \in \mathcal{E}_T} |\lambda_{KL}^T| |\xi(u_K^{n+1}) - \xi(u_L^{n+1})|^2 \\ &\leq C \mathcal{Q}_h \longrightarrow 0 \quad \text{as } h \rightarrow 0. \end{aligned}$$

Convection term For all $T \in \mathcal{T}_h$, we define the piecewise constant function $\kappa_{\mathcal{T}_h, \Delta t_h}$ by

$$\kappa_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t) = \chi \circ \xi^{-1}(\Upsilon_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}, t)), \quad \forall \mathbf{x} \in T, \forall t \in (t_n, t_{n+1}].$$

Using the same guidelines as for the convergence results (4.64), one has

$$\kappa_{\mathcal{T}_h, \Delta t_h} \longrightarrow \chi(u) \quad \text{in } L^2(Q_{t_f}) \text{ as } h \rightarrow 0.$$

We introduce the term

$$\mathcal{D}'_h := - \iint_{Q_{t_f}} \kappa_{\mathcal{T}_h, \Delta t_h} \Lambda(\mathbf{x}) \nabla v_{\mathcal{T}_h, \Delta t_h} \cdot \nabla \psi_{\mathcal{T}_h, \Delta t_h}(\cdot, t - \Delta t_h) d\mathbf{x} dt.$$

Thanks to the weakly convergence in $L^2(Q_{t_f})$ of the sequence $\nabla v_{\mathcal{T}_h, \Delta t_h}$ towards ∇v , and to the uniform convergence of $\nabla \psi_{\mathcal{T}_h, \Delta t_h}$ towards $\nabla \psi$, we obtain

$$\mathcal{D}'_h \longrightarrow - \iint_{Q_{t_f}} \chi(u) \Lambda(\mathbf{x}) \nabla v \cdot \nabla \psi d\mathbf{x} dt \quad \text{as } h \rightarrow 0.$$

Let us prove, using the same guidelines as before, that $|\mathcal{D}_h - \mathcal{D}'_h| \longrightarrow 0$, when h tends to zero.

We denote by

$$\begin{aligned} \chi_T^{n+1} &= \kappa_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}_T, t_{n+1}), & \forall T \in \mathcal{T}_h, \forall n \in \{0, \dots, N_h\}, \\ \mu_T^{n+1} &= \mu_{\mathcal{T}_h, \Delta t_h}(\mathbf{x}_T, t_{n+1}), & \forall T \in \mathcal{T}_h, \forall n \in \{0, \dots, N_h\}. \end{aligned}$$

Therefore,

$$\mathcal{D}_h - \mathcal{D}'_h = \sum_{n=0}^{N_h} \Delta t_h \sum_{T \in \mathcal{T}_h} \sum_{\sigma_{KL} \in \mathcal{E}_T} (a_T^{n+1} \mu_T^{n+1} - a_{KL}^{n+1} \mu_{KL}^{n+1}) \lambda_{KL}^T (v_K^{n+1} - v_L^{n+1}) (\psi_K^n - \psi_L^n).$$

Thanks to the triangle inequality and to the existence of a continuity moduli η and δ of the continuous functions $\sqrt{a \circ \xi^{-1}}$ and $\mu \circ \xi^{-1}$ respectively, one has

$$\begin{aligned} |a_{KL}^{n+1} \mu_{KL}^{n+1} - a_T^{n+1} \mu_T^{n+1}| &\leq \mu_{KL}^{n+1} |a_{KL}^{n+1} - a_T^{n+1}| + a_T^{n+1} |\mu_{KL}^{n+1} - \mu_T^{n+1}| \\ &\leq C \left(\eta \left(\bar{\xi}_T^{n+1} - \underline{\xi}_T^{n+1} \right) + \delta \left(\bar{\xi}_T^{n+1} - \underline{\xi}_T^{n+1} \right) \right), \end{aligned}$$

where the constant C does not depend on h . Therefore, using the Cauchy-Schwarz inequality, Lemma 4.3, Lemma 4.4, and Proposition 4.9, we deduce that there exists a constant C independent of h such that

$$|\mathcal{D}_h - \mathcal{D}'_h|^2 \leq C(Q_h + \mathcal{W}_h) \times \sum_{n=0}^{N_h} \Delta t_h \sum_{T \in \mathcal{T}_h} \sum_{\sigma_{KL} \in \mathcal{E}_T} |\lambda_{KL}^T| |v_K^{n+1} - v_L^{n+1}|^2,$$

where Q_h is given by equation (4.63), and \mathcal{W}_h is given by

$$\mathcal{W}_h = \sum_{n=0}^{N_h} \Delta t_h \sum_{T \in \mathcal{T}_h} \left(\delta \left(\bar{\xi}_T^{n+1} - \underline{\xi}_T^{n+1} \right) \right)^2 \sum_{\sigma_{KL} \in \mathcal{E}_T} |\lambda_{KL}^T| (\psi_K^n - \psi_L^n)^2.$$

Now, using the same proof as for the diffusive term, one can deduce that $\mathcal{W}_h \leq Ch$. Therefore

$$\lim_{h \rightarrow 0} |\mathcal{D}_h - \mathcal{D}'_h| = 0,$$

and consequently,

$$\lim_{h \rightarrow 0} \mathcal{D}_h = - \iint_{Q_{t_f}} \chi(u) \Lambda(\mathbf{x}) \nabla v \cdot \nabla \psi d\mathbf{x} dt.$$

Reaction term We would now like to show that

$$\mathcal{F}_h \longrightarrow \iint_{Q_{t_f}} f(u(\mathbf{x}, t)) \psi(\mathbf{x}, t) \, d\mathbf{x} \, dt \quad \text{as } h \rightarrow 0.$$

For this purpose, we denote, for all $K \in \mathcal{V}_h$ and for all $n \geq 1$, by $f_K^n = f(u_K^n)$, and by $f_{\mathcal{M}_h, \Delta t_h}$ the piecewise constant reconstruction in $\mathcal{X}_{\mathcal{M}_h, \Delta t_h}$. Thus we have

$$\mathcal{F}_h = \iint_{Q_{t_f}} f_{\mathcal{M}_h, \Delta t_h} \psi_{\mathcal{M}_h, \Delta t_h}(\cdot, t - \Delta t_h) \, d\mathbf{x} \, dt \longrightarrow \iint_{Q_{t_f}} f(u(\mathbf{x}, t)) \psi(\mathbf{x}, t) \, d\mathbf{x} \, dt \quad \text{as } h \rightarrow 0,$$

since $f(u)_{\mathcal{M}_h, \Delta t_h}$ converges strongly in $L^2(Q_{t_f})$ towards $f(u)$, and as $\psi_{\mathcal{M}_h, \Delta t_h}$ converges uniformly towards ψ . This ends the proof of the convergence.

4.8 Numerical results

In this section, we establish various 2-D numerical results provided by the *nonlinear CVFE scheme* (4.13)–(4.21). Newton's algorithm is carried out for the implementation of the scheme, coupled with a biconjugate gradient method to solve linear systems arising from the Newton algorithm. We provide three tests to show the effectiveness of the *nonlinear CVFE scheme* (4.13)–(4.21). For these tests, we consider the following data : $L_x = 1$, $L_y = 1$ (the length and the width of the domain). We fix : $\Delta t = 0.002$, $\alpha = 0.01$, $\beta = 0.05$, $a(u) = d_u u(1 - u)$, $d_u = 0.0005$, $\chi(u) = \zeta \times (u(1 - u))^2$, $\zeta = 0.05$. By definition, we have $\mu(u) = \frac{\zeta}{d_u} u(1 - u)$ then, the numerical flux function μ_{KL}^{n+1} is given using the following functions :

$$\mu_{\uparrow}(z) = \mu\left(\min\left\{z, \frac{1}{2}\right\}\right), \quad \text{and} \quad \mu_{\downarrow}(z) = \mu\left(\max\left\{z, \frac{1}{2}\right\}\right) - \mu\left(\frac{1}{2}\right), \quad \forall z \in (0, 1) \times (0, 1).$$

Unless stated otherwise and throughout the tests, we assume that $f(u) = 0$, that the initial conditions are defined by regions, and we assume zeros-flux boundary conditions. For instance, the cell density is initially defined by $u_0(\mathbf{x}, \mathbf{y}) = 1$ in the square region given by $(\mathbf{x}, \mathbf{y}) \in [0.45, 0.55]$ and 0 otherwise. The initial chemoattractant concentration is defined by $v_0(\mathbf{x}, \mathbf{y}) = 5$ in the space region given by $(\mathbf{x}, \mathbf{y}) \in [0.2, 0.3] \times [0.45, 0.55] \cup [0.45, 0.55] \times [0.2, 0.3] \cup [0.45, 0.55] \times [0.7, 0.8] \cup [0.7, 0.8] \times [0.45, 0.55]$.

Test 1 (Weak anisotropic case). In this test, we assume that the diffusion tensors are given by

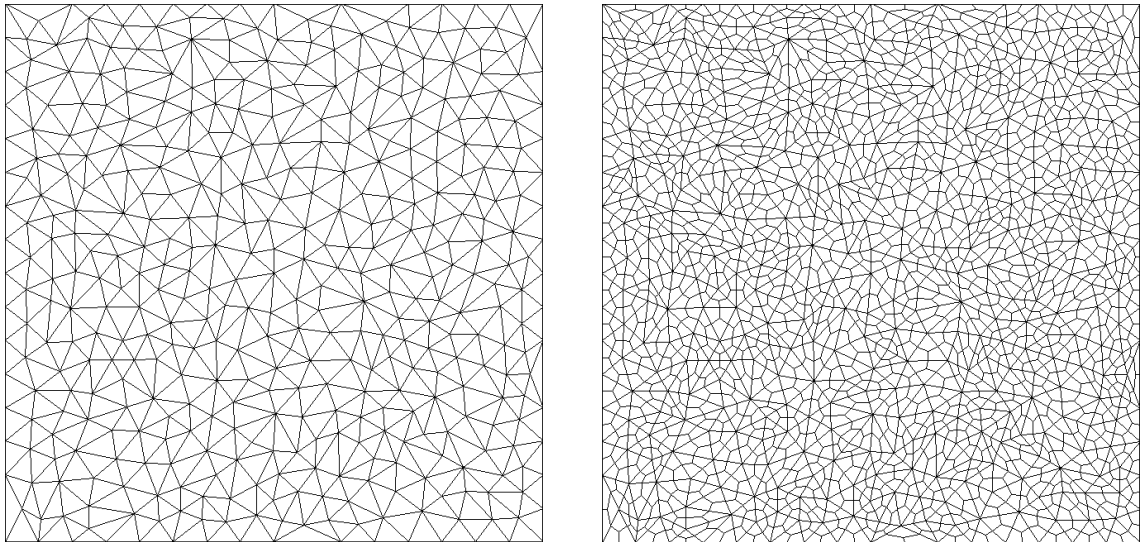
$$\Lambda(\mathbf{x}) = \begin{pmatrix} 1 & 0 \\ 0 & \theta \end{pmatrix}, \quad D(\mathbf{x}) = dI_2, \quad d = 0.0001.$$

Further, we consider an admissible triangular primary mesh made of 14 336 triangles, the corresponding Donald dual mesh consists of 7 297 dual control volumes. In a admissible triangular mesh, all the angles of triangles are acute, then one can deduce that the transmissibility coefficients are nonnegative in the case of isotropic diffusion tensors. In Tab. 4.1, we present minimum and maximum values obtained with each of the scheme (3.24)–(3.25) (of the previous chapter), the *nonlinear CVFE scheme* (4.13)–(4.21), and the FV scheme.

Test 2 (Weak anisotropic case/obtuse angles). In this test, we consider a general unstructured mesh that contains obtuse angles, this mesh is made of 5 193 triangles and 2 665 dual control volumes (see Figure 4.2 for an illustration of the primal and dual mesh).

		<i>scheme (3.24)–(3.25)</i>	<i>scheme (4.13)–(4.21)</i>	FV scheme
After 1 iteration $\theta = 1$	Min. Val. u Max. Val. u	0.000000 1.000000	0.000000 1.000000	0.000000 1.000000
After 10 iterations $\theta = 1$	Min. Val. u Max. Val. u	0.000000 0.971110	0.000000 1.000000	0.000000 1.000000
After 1 iteration $\theta = 5$	Min. Val. u Max. Val. u	-1.73001×10^{-3} 0.99722922	8.68789×10^{-20} 1.000000	
After 10 iterations $\theta = 5$	Min. Val. u Max. Val. u	-1.62500×10^{-2} 0.9715705	0.000000 1.000000	
After 1 iteration $\theta = 10$	Min. Val. u Max. Val. u	-4.46953×10^{-3} 1.00018368	6.30555×10^{-16} 1.000000	
After 10 iterations $\theta = 10$	Min. Val. u Max. Val. u	-3.91245×10^{-2} 0.98342428	6.30554×10^{-16} 0.9999999	

TABLE 4.1 – Numerical results after 1 and 10 iterations.

FIGURE 4.2 – Initial primal mesh \mathcal{T}_h and barycentric dual mesh \mathcal{M}_h .

The diffusion tensors are defined, for all $\mathbf{x} \in (0, 1) \times (0, 1)$, by

$$\Lambda(\mathbf{x}) = \begin{pmatrix} 7 & 2 \\ 2 & 10 \end{pmatrix}, \quad D(\mathbf{x}) = dI_2, \quad d = 0.0001.$$

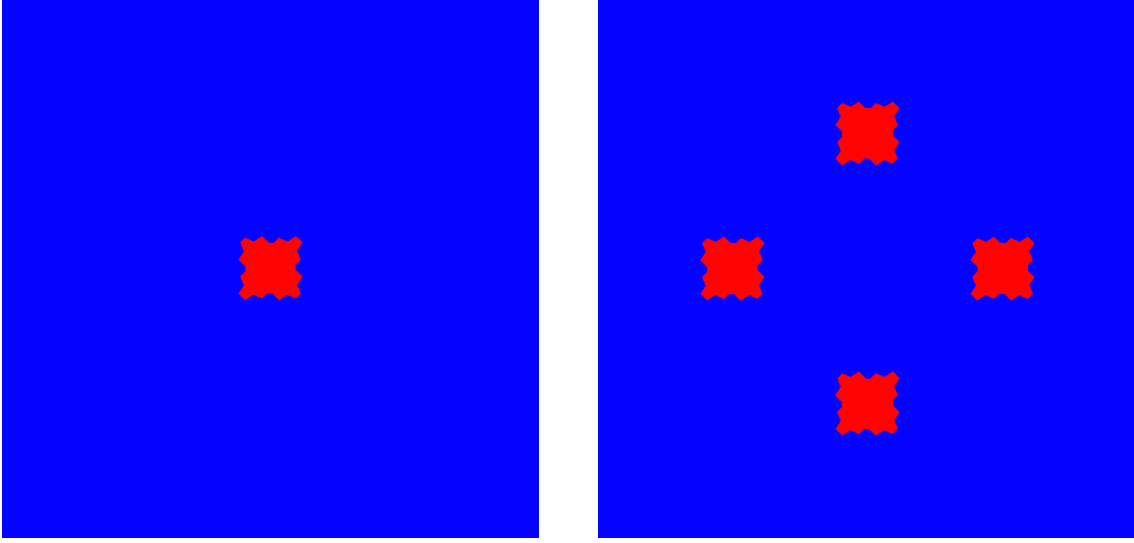


FIGURE 4.3 – Initial condition for the cell density u (left) with $0 \leq u \leq 1$ and for the chemoattractant concentration v (right) with $0 \leq v \leq 5$.

The cell density diffusion tensor Λ is taken to be a homogeneous anisotropic tensor with high diffusivity in a direction at 63.44 degrees from the horizontal and low diffusivity in the orthogonal direction. Indeed, we have

$$\Lambda = R_\theta \times \begin{pmatrix} 11 & 0 \\ 0 & 6 \end{pmatrix} \times R_\theta^{-1},$$

where, $R_\theta = \begin{pmatrix} 0.4472 & -0.8944 \\ 0.8944 & 0.4472 \end{pmatrix} = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}$ is the rotation of angle $\theta = 63.44$ degrees.

Figure 4.3 represents initial distributions of the cell density u and the chemoattractant concentration v over the initial triangular mesh as well as the corresponding dual mesh.

Figures 4.4–4.5 represent the evolution of the cell density at time $t = 0.4$, $t = 1.4$, $t = 2.4$, and $t = 4$. At moment $t = 0.4$, it is clear that the cell density diffuses in the space (in a direction at 63.44 degrees from the horizontal) without any interactions with the chemoattractant which diffuses uniformly in the space. Then, after a while, and when the chemoattractant diffusion reaches the cell density location, we see that the latter changes its direction to be absorbed by the chemoattractant located vertically. This process continues and the cells accumulate into the location of the chemoattractant and finally we obtain the cell density aggregations as shown at $t = 4$.

Test 3 (Anisotropic case/obtuse angles). In this test, we consider an unstructured mesh consisting of 15 568 primal triangles and 7 912 dual control dual volumes. Further, we assume that the diffusion tensors are anisotropic and are given by :

$$\Lambda(\mathbf{x}) = \begin{pmatrix} 8 & -7 \\ -7 & 20 \end{pmatrix}, \quad D(\mathbf{x}) = d \begin{pmatrix} 1 & 0 \\ 0 & \theta \end{pmatrix}, \quad d = 0.0001.$$

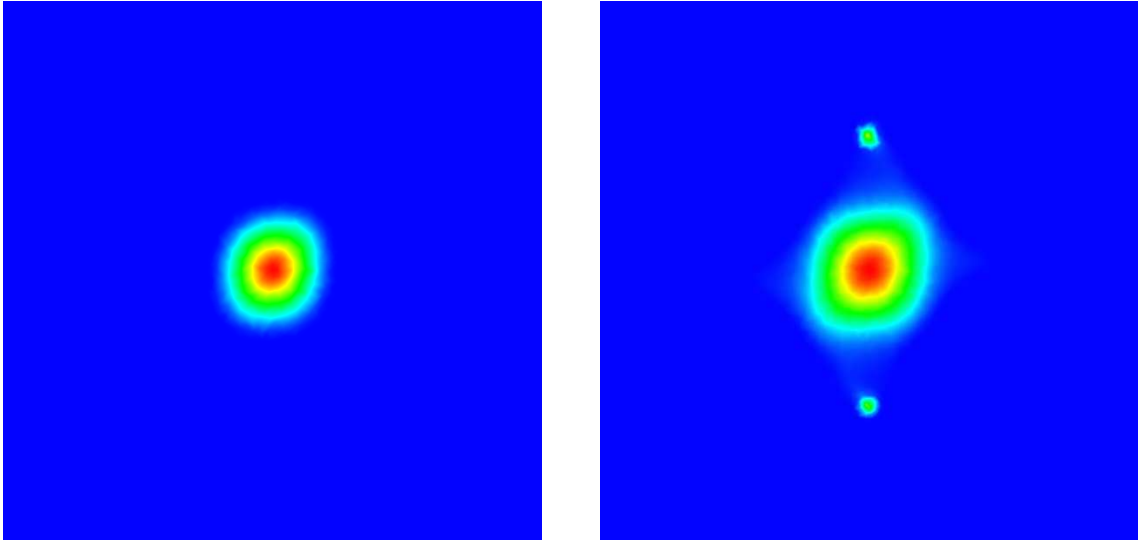


FIGURE 4.4 – Evolution of the cell density u at time $t = 0.4$ with $0 \leq u \leq 0.667$ (left), and at time $t = 1.4$ with $0 \leq u \leq 0.632$ (right).

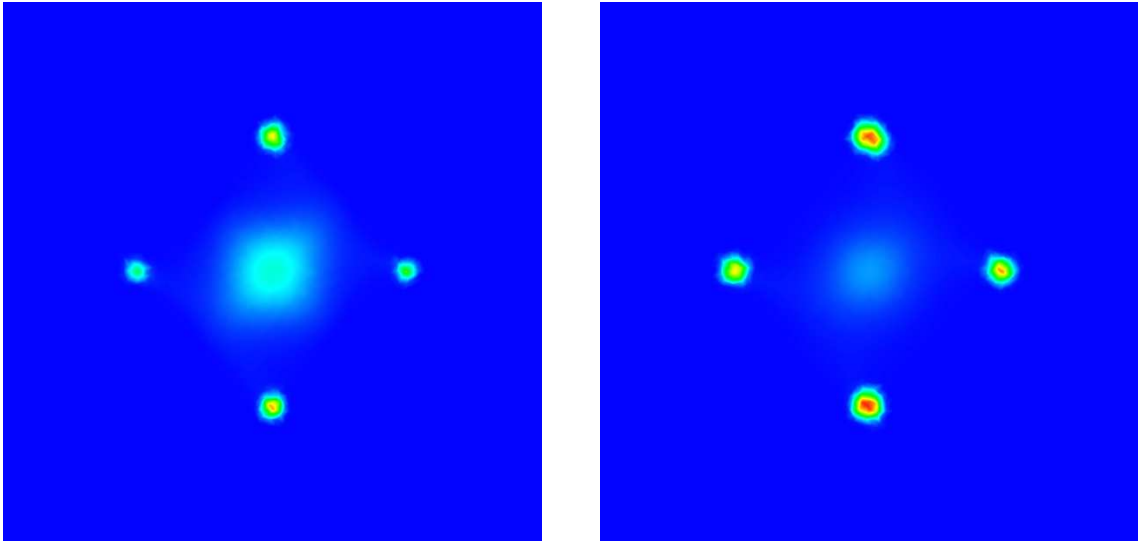


FIGURE 4.5 – Evolution of the cell density u at time $t = 2.4$ with $0 \leq u \leq 0.972$ (left), and at time $t = 4$ with $0 \leq u \leq 0.987$ (right).

Table 4.2 provides a comparison between the *nonlinear CVFE scheme* coupled on the one hand with the discretization (4.15) of v and with the discretization (4.21) of v on the other hand. We see that the discretization (4.21) carries out a better approximation than the discretization (4.15) in terms of ensuring the discrete maximum principle property.

		<i>scheme (4.13)–(4.15)</i>	<i>scheme (4.13)–(4.21)</i>
After 1 iteration $\theta = 3$	Min. Val. u	0.0000000	0.0000000
	Max. Val. u	1.0000000	1.0000000
	Min. Val. v	-1.141912E-002	1.922764E-051
	Max. Val. v	5.012383	4.999982
After 200 iterations $\theta = 3$	Min. Val. u	0.0000000	0.0000000
	Max. Val. u	0.5298226	0.5312562
	Min. Val. v	-1.731068E-003	1.297192E-080
	Max. Val. v	4.8053827	4.8018742
After 1000 iterations $\theta = 3$	Min. Val. u	0.0000000	0.0000000
	Max. Val. u	0.9957580	0.9974757
	Min. Val. v	6.265859E-023	3.171769E-080
	Max. Val. v	2.961761	2.910828

TABLE 4.2 – Numerical results after 1, 200 and 1000 iterations over an unstructured mesh with obtuse angles.

Analysis of a nonlinear degenerate parabolic equation arising in chemotaxis or in porous media

Sommaire

5.1	Introduction	119
5.2	The nonlinear degenerate model	120
5.2.1	Classical weak solutions	121
5.2.2	Weak degenerate solutions	122
5.3	Existence for the nondegenerate case	123
5.3.1	Weak nondegenerate solutions	124
5.3.2	Maximum principle on the saturation	128
5.4	Proof of theorem 5.2.	130
5.5	Proof of theorem 5.4	134
5.6	Proof of theorem 5.6	145

5.1 Introduction

This chapter is devoted to the theoretical analysis of a general degenerate nonlinear parabolic equation. This kind of equations stems either from the modeling of a compressible two phase flow in porous media or from the modeling of the chemotaxis-fluid process.

In the degenerate equation, the strong nonlinearities are technically difficult to be controlled by the degenerate dissipative term because the equation itself presents degenerate terms of order 0 and of order 1.

In the case of the degeneracy of the dissipative term at one point, a weak and classical formulation is possible for the expected solutions. However, in the case of the degeneracy of the dissipative term at two points, we obtain solutions in a weaker sense compared to the one of the

classical formulation. Therefore, a degenerate weighted formulation is introduced taking into account the degeneracy of the dissipative term.

5.2 The nonlinear degenerate model

Let $T > 0$ be a fixed time and Ω be an open bounded subset of \mathbb{R}^d , $d = 2, 3$. We set $Q_T := \Omega \times (0, T)$ and $\Sigma_T = \partial\Omega \times (0, T)$. We consider the following nonlinear degenerate parabolic equation

$$\partial_t u - \operatorname{div} (a(u) \nabla u - f(u) \mathbf{V}) - g(u) \operatorname{div} (\mathbf{V}) + \gamma a(u) \nabla u \cdot \tilde{\mathbf{V}} = 0, \quad \text{in } Q_T. \quad (5.1)$$

The boundary condition is defined, on the saturation u , by

$$u(\mathbf{x}, t) = 0, \quad \text{on } \Sigma_T. \quad (5.2)$$

The initial condition is given by

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad \text{in } \Omega. \quad (5.3)$$

Many physical models lead to degenerate nonlinear parabolic problem. For instance, in [39] the authors analyzed a model of a degenerate nonlinear system arising from compressible two-phase flows in porous media (in this case we have $\mathbf{V} = \tilde{\mathbf{V}}$). The described system coupled the saturation (denoted by u) and the global pressure (denoted by p). The global velocity (denoted by \mathbf{V}) is taken to be proportional to the gradient of the global pressure. On the other hand, several papers are devoted to the mathematical analysis of compressible, miscible displacement models in porous media (see e.g. [5, 6, 35]). Here, we consider a generalization of the saturation equation where we assume that the velocity field is given and fixed.

Other models that can lead to such kind of degenerate nonlinear parabolic equation (5.1) are the chemotaxis models, where u represents the cell density and \mathbf{V} represents the gradient of the chemical concentration (see e.g. [8, 26, 57]), while $\tilde{\mathbf{V}}$ represents the velocity of the fluid which transports the cell density and the chemical concentration.

We introduce the classical assumptions for porous media made about the system (5.1)–(5.3) :

(H1) The velocities \mathbf{V} and $\tilde{\mathbf{V}}$ are two measurable functions lying into $(L^\infty(\Omega))^d$. In addition, we assume the following condition on the velocity field \mathbf{V} :

$$\mathbf{V} \cdot \mathbf{n} \leq 0, \quad \text{on } \Sigma_T,$$

where \mathbf{n} is the unit normal vector to $\partial\Omega$ outward to Ω .

(H2) f is a differentiable function in $[0, 1]$ and $g \in C^1([0, 1])$ verifying

$$g(0) = f(0) = 0, \quad f(1) = g(1) = 1, \quad \text{and } g'(u) \geq C_{g'} > 0 \quad \forall u \in [0, 1].$$

(H3) The initial condition u_0 satisfies : $u_0 \in L^2(\Omega)$ and $0 \leq u_0(\mathbf{x}) \leq 1$ for a.e. $\mathbf{x} \in \Omega$.

A major difficulty of system (5.1)–(5.3) is the possible degeneracy of the diffusion term. Here, we give the degeneracy assumption on the dissipation function a :

(H4a) $a \in C^1([0, 1], \mathbb{R})$, $a(u) > 0$ for $0 < u < 1$, $a(0) > 0$, $a(1) = 0$.

Furthermore, there exist $a_0 > 0$, $0 < r_2 \leq 2$, $u_* < 1$, m_1 and $M_1 > 0$ such that

$$a(u) \geq a_0 \quad \text{for all } 0 \leq u \leq u_*,$$

$$m_1 (1 - u)^{r_2} \leq a(u) \leq M_1 (1 - u)^{r_2}, \quad \text{for all } u_* \leq u \leq 1. \quad \text{In addition there exists } c_1, c_2 > 0 \text{ such that } c_1 (1 - u)^{-1} \leq (f(u) - g(u))^{-1} \leq c_2 (1 - u)^{-1} \text{ for all } u_* \leq u < 1.$$

(H4b) $a \in \mathcal{C}^1([0, 1], \mathbb{R})$, $a(u) > 0$ for $0 < u < 1$, $a(0) = 0$, $a(1) > 0$,

Furthermore, there exist $r_1 > 0$, m_1 and $M_1 > 0$ such that

$m_1 r_1 u^{r_1-1} \leq a'(u) \leq M_1 r_1 u^{r_1-1}$, for all $0 \leq u \leq 1$. In addition, there exists a constant $C > 0$ such that $|f(u) - g(u)| \leq Cu$, $\forall 0 \leq u \leq 1$.

(H4c) $a \in \mathcal{C}^1([0, 1], \mathbb{R})$, $a(u) > 0$ for $0 < u < 1$, $a(0) = 0$, $a(1) = 0$,

Furthermore, there exist $r_1 > 0$, $r_2 > 0$, $u_* < 1$, m_1 and $M_1 > 0$ such that

$m_1 r_1 u^{r_1-1} \leq a'(u) \leq M_1 r_1 u^{r_1-1}$, for all $0 \leq u \leq u_*$,

$-r_2 M_1 (1-u)^{r_2-1} \leq a'(u) \leq -r_2 m_1 (1-u)^{r_2-1}$, for all $u_* \leq u \leq 1$. In addition, there exists $c_1, c_2, C > 0$ such that $|f(u) - g(u)| \leq Cu$ for all $0 \leq u \leq u_*$ et $c_1 (1-u)^{-1} \leq (f(u) - g(u))^{-1} \leq c_2 (1-u)^{-1}$ for all $u_* \leq u < 1$.

In what follows, we introduce first the existence of classical weak solutions to equation (5.1) under the assumptions (H1)–(H3) and (H4a) and for a particular choice of the initial data. Next, we show the existence of solutions to equation (5.1) in a weak sense (by introducing a weighted formulation) and under the assumptions (H1)–(H3) and (H4b). Finally, for a particular choice of the initial data, we introduce also the existence of solutions to equation (5.1) (verifying a weighted formulation) and under the assumptions (H1)–(H3) and (H4c).

In the sequel and for the simplicity, we assume that the nonnegative constant γ is fixed equals to 1 and that $\tilde{\mathbf{V}} = \mathbf{V}$ (the same analysis is possible with $\tilde{\mathbf{V}} \neq \mathbf{V}$).

5.2.1 Classical weak solutions

Let us consider the function $k \in \mathcal{C}^1[0, 1)$ defined by

$$\begin{aligned} k(u) &= u, & \text{if } 0 \leq u \leq u_*, \\ k'(u) &= (f(u) - g(u))^{-1} g'(u) k(u), & \text{if } u_* \leq u < 1. \end{aligned} \quad (5.4)$$

We denote by L the following primitive of the function k

$$L(u) = \int_0^u k(\tau) \, d\tau, \quad \forall 0 \leq u < 1.$$

Definition 5.1. Under the assumptions (H1)–(H3) and (H4a), and assume that the initial condition $u(t=0) = u_0$ satisfies $L(u_0) \in L^1(\Omega)$. We say that u is a classical weak solution to system (5.1)–(5.3) if u verifies

$$\begin{aligned} 0 &\leq u(\mathbf{x}, t) \leq 1 \text{ for a.e. } (\mathbf{x}, t) \in \Omega \times (0, T), \\ u &\in L^2(0, T; H_0^1(\Omega)) \cap \mathcal{C}^0([0, T]; L^2(\Omega)), \\ \partial_t u &\in L^2(0, T; H^{-1}(\Omega)), \end{aligned}$$

and such that

$$\begin{aligned} \int_0^T \langle \partial_t u, \varphi \rangle \, dt &+ \int_{Q_T} a(u) \nabla u \cdot \nabla \varphi \, d\mathbf{x} \, dt - \int_{Q_T} f(u) \mathbf{V} \cdot \nabla \varphi \, d\mathbf{x} \, dt \\ &+ \int_{Q_T} g(u) \mathbf{V} \cdot \nabla \varphi \, d\mathbf{x} \, dt + \int_{Q_T} g'(u) \mathbf{V} \cdot \nabla u \varphi \, d\mathbf{x} \, dt \\ &+ \int_{Q_T} a(u) \mathbf{V} \cdot \nabla u \varphi \, d\mathbf{x} \, dt = 0, \quad \forall \varphi \in L^2(0, T; H_0^1(\Omega)). \end{aligned} \quad (5.5)$$

where the bracket $\langle \cdot, \cdot \rangle$ represents the duality product between $H^{-1}(\Omega)$ and $H_0^1(\Omega)$.

Theorem 5.2 (Degenerate system). *Under the assumptions (H1) – (H3) and (H4a), there exists at least one classical weak solution to the degenerate system (5.1)–(5.3) in the sense of Definition 5.1.*

5.2.2 Weak degenerate solutions

Let us introduce the functions β and h defined by

$$\beta(u) = u^{r-1}, \quad h(u) = \int_0^u \beta(\tau) d\tau, \quad \text{where } r = \begin{cases} r_1 + 2, & \text{if } r_1 \leq 1, \\ r_1, & \text{if } r_1 > 1. \end{cases}$$

Here, r_1 is the same constant defined in assumption (H4b).

Further, we consider the functions β_θ and h_θ defined, for all $\theta > 0$, by

$$\beta_\theta(u) = u^{r-1+\theta}, \quad h_\theta(u) = \int_0^u \beta_\theta(\tau) d\tau.$$

Definition 5.3. For $\theta \geq 7r_1 + 6 - r$, and under the assumptions (H1)–(H3) and (H4b). We say that u is a degenerate weak solution to system (5.1)–(5.3) if

$$\begin{aligned} 0 &\leq u(\mathbf{x}, t) \leq 1 \text{ for a.e. } (\mathbf{x}, t) \in \Omega \times (0, T), \\ h_\theta(u) &\in L^2(0, T; H_0^1(\Omega)), \\ \sqrt{\beta'(u)} \nabla u &\in (L^2(Q_T))^d, \end{aligned}$$

and such that the function F defined, for all $\chi \in L^2(0, T; H^1(\Omega))$ by

$$\begin{aligned} F(u, \chi) = & - \int_{Q_T} h_\theta(u) \partial_t \chi d\mathbf{x} dt - \int_\Omega h_\theta(u_0) \chi(\mathbf{x}, 0) d\mathbf{x} + \int_{Q_T} a(u) \nabla u \cdot \nabla (\beta_\theta(u) \chi) d\mathbf{x} dt \\ & + \int_{Q_T} a(u) \mathbf{V} \cdot \nabla u \beta_\theta(u) \chi d\mathbf{x} dt - \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla (\beta_\theta(u) \chi) d\mathbf{x} dt \\ & + \int_{Q_T} g'(u) \mathbf{V} \cdot \nabla u \beta_\theta(u) \chi d\mathbf{x} dt - \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla \chi \beta_\theta(u) d\mathbf{x} dt, \end{aligned}$$

verifies

$$F(u, \chi) \leq 0, \quad \forall \chi \in \mathcal{C}^1([0, T]; H_0^1(\Omega)) \text{ with } \chi(\cdot, T) = 0 \text{ and } \chi \geq 0. \quad (5.6)$$

Furthermore, it is required to satisfy

$$\begin{aligned} \forall \varepsilon > 0, \exists Q^\varepsilon \subset Q_T \text{ such that } \text{meas}(Q^\varepsilon) < \varepsilon, \text{ and} \\ F(u, \chi) = 0, \quad \forall \chi \in \mathcal{C}^1([0, T]; H_0^1(\Omega)), \text{ supp } \chi \subset ([0, T] \times \Omega) \setminus Q^\varepsilon \end{aligned} \quad (5.7)$$

Theorem 5.4. *Under the assumptions (H1) – (H3) and (H4b), there exists at least one weak degenerate solution to system (5.1)–(5.3) in the sense of Definition 5.3.*

We give now the definition of weak solutions to system (5.1)–(5.3) when the assumption (H4c) is satisfied. Let $\theta, \lambda \geq 0$, we denote by $j_{\theta, \lambda}$ the continuous function defined by

$$j_{\theta, \lambda}(u) = \begin{cases} \beta_\theta(u), & \text{if } 0 \leq u \leq u_* \\ \beta_\theta(u_*) (1 - u_*)^{1 - \frac{r'}{2} - \lambda} (1 - u)^{\frac{r'}{2} - 1} + \lambda, & \text{if } u \geq u_*. \end{cases} \quad (5.8)$$

where $r' \geq \max(2, r_2)$. We denote by $J_{\theta, \lambda}$ the primitive of the function $j_{\theta, \lambda}$

$$J_{\theta, \lambda} = \int_0^u j_{\theta, \lambda}(y) dy. \quad (5.9)$$

To simplify the notations, we denote by $j = j_{0,0}$ and $J = J_{0,0}$. In addition, we consider the functions μ and G , defined by

$$\begin{cases} \mu(u) = \beta(u), & 0 \leq u \leq u_* \\ \mu'(u) = (f(u) - g(u))^{-1} g'(u) \mu(u), & u_* \leq u < 1. \end{cases} \quad (5.10)$$

G is the primitive of μ , that is

$$G(u) = \int_0^u \mu(y) \, dy. \quad (5.11)$$

Definition 5.5. For $\theta \geq 7r_1 + 6 - r$, $\lambda \geq 7r_2 + 6 - \frac{r'}{2}$, and under the assumptions (H1)–(H3) and (H4c). We say that u is a degenerate weak solution to system (5.1)–(5.3) if

$$\begin{aligned} 0 \leq u(\mathbf{x}, t) \leq 1 \text{ for a.e. } (\mathbf{x}, t) \in \Omega \times (0, T), \\ J(u) \in L^2(0, T; H_0^1(\Omega)), \quad \mu'^{\frac{1}{2}}(u) a^{\frac{1}{2}}(u) \nabla u \in (L^2(Q_T))^d, \end{aligned}$$

and such that, the function F defined by

$$\begin{aligned} F(u, \chi) = & - \int_{Q_T} J_{\theta, \lambda}(u) \partial_t \chi \, d\mathbf{x} \, dt - \int_{\Omega} J_{\theta, \lambda}(u_0(\mathbf{x})) \chi(\mathbf{x}, 0) \, d\mathbf{x} \\ & + \int_{Q_T} a(u) \nabla u \cdot \nabla (j_{\theta, \lambda}(u) \chi) \, d\mathbf{x} \, dt + \int_{Q_T} a(u) \mathbf{V} \cdot \nabla u j_{\theta, \lambda}(u) \chi \, d\mathbf{x} \, dt \\ & - \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla (j_{\theta, \lambda}(u) \chi) \, d\mathbf{x} \, dt + \int_{Q_T} g'(u) \mathbf{V} \cdot \nabla u j_{\theta, \lambda}(u) \chi \, d\mathbf{x} \, dt \\ & - \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla \chi j_{\theta, \lambda}(u) \, d\mathbf{x} \, dt, \end{aligned}$$

verifies

$$F(u, \chi) \leq 0, \quad \forall \chi \in \mathcal{C}^1([0, T]; H_0^1(\Omega)) \text{ with } \chi(\cdot, T) = 0 \text{ and } \chi \geq 0 \quad (5.12)$$

and furthermore,

$$\begin{aligned} \forall \varepsilon > 0, \exists Q^\varepsilon \subset Q_T \text{ such that } \text{meas}(Q^\varepsilon) < \varepsilon, \text{ and} \\ F(u, \chi) = 0, \quad \forall \chi \in \mathcal{C}^1([0, T]; H_0^1(\Omega)), \text{ supp } \chi \subset ([0, T] \times \Omega) \setminus Q^\varepsilon \end{aligned} \quad (5.13)$$

Theorem 5.6. Under the assumptions (H1) – (H3) and (H4c), there exists at least one degenerate weak solution to system (5.1)–(5.3) in the sense of Definition 5.5.

Unless stated otherwise, κ represents a “generic” nonnegative quantity which need not have the same value through the proofs. Furthermore, C_α represents a nonnegative constant depending only on the subscript α .

5.3 Existence for the nondegenerate case

In this section, we prove the existence of solutions to the nondegenerate problem. To avoid the degeneracy of the dissipation function a , we introduce the modified problem where the dissipation a is replaced by $a_\eta(u) = a(u) + \eta$ in equation (5.1), with $0 < \eta \ll 1$.

Therefore, we consider the nondegenerate system

$$\partial_t u_\eta - \text{div}(a_\eta(u_\eta) \nabla u_\eta - f(u_\eta) \mathbf{V}) - g(u_\eta) \text{div}(\mathbf{V}) + a_\eta(u_\eta) \nabla u_\eta \cdot \mathbf{V} = 0, \text{ in } Q_T, \quad (5.14)$$

$$u_\eta(\mathbf{x}, t) = 0, \quad \mathbf{V} \cdot \mathbf{n} \leq 0, \quad \text{in } \Sigma_T, \quad (5.15)$$

$$u_\eta(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad \text{in } \Omega. \quad (5.16)$$

5.3.1 Weak nondegenerate solutions

For the existence of a solution to the nondegenerate system, we make the following weak assumption

(H5) $a \in \mathcal{C}^0([0, 1])$, $a(u) > 0$ for $0 < u < 1$, $a(0) = 0$, and $a(1) = 0$ or $a(1) > 0$.

Theorem 5.7 (nondegenerate system). *For any fixed $\eta > 0$ and under the assumptions (H1)–(H3) and (H5), there exists at least one weak solution u_η to the system (5.14)–(5.16) satisfying*

$$\begin{aligned} 0 &\leq u_\eta(\mathbf{x}, t) \leq 1 \quad \text{for a.e. } (\mathbf{x}, t) \in Q_T, \\ u_\eta &\in L^2(0, T; H_0^1(\Omega)) \cap L^\infty(0, T; L^2(\Omega)), \\ \partial_t u_\eta &\in L^2(0, T; H^{-1}(\Omega)), \\ u_\eta &\in \mathcal{C}^0([0, T]; L^2(\Omega)), \end{aligned} \quad (5.17)$$

and such that

$$\begin{aligned} \int_0^T \langle \partial_t u_\eta, \varphi \rangle d\tau + \int_{Q_T} a_\eta(u_\eta) \nabla u_\eta \cdot \nabla \varphi d\mathbf{x} dt - \int_{Q_T} f(u_\eta) \mathbf{V} \cdot \nabla \varphi d\mathbf{x} dt \\ + \int_{Q_T} g'(u_\eta) \mathbf{V} \cdot \nabla u_\eta \varphi d\mathbf{x} dt + \int_{Q_T} g(u_\eta) \mathbf{V} \cdot \nabla \varphi d\mathbf{x} dt \\ + \int_{Q_T} a_\eta(u_\eta) \mathbf{V} \cdot \nabla u_\eta \varphi d\mathbf{x} dt = 0, \quad \forall \varphi \in L^2(0, T; H_0^1(\Omega)), \end{aligned} \quad (5.18)$$

where the bracket $\langle \cdot, \cdot \rangle$ represents the duality product between $H^{-1}(\Omega)$ and $H_0^1(\Omega)$.

Proof. The solutions to system (5.14)–(5.16) depend on the parameter η . To simplify the notations and for simplicity, we omit the dependence of solutions on the parameter η and we use u instead of u_η in this section. We will apply the Schauder fixed-point theorem to prove the existence of weak solutions to system (5.14)–(5.16).

It is necessary to use the continuous extension for the functions depending on u . For instance, we take $f(u) = g(u) = 1$ for all $u \geq 1$ and $f(u) = g(u) = 0$ for all $u \leq 0$. Furthermore, we extend the dissipation a outside $[0, 1]$ by taking

$$a(u) = 0, \text{ for } u \leq 0, \quad \text{and} \quad a(u) = a(1), \text{ for } u \geq 0.$$

Fixed-point method

Let us introduce the closed subset \mathcal{K} of $L^2(Q_T)$ given by

$$\mathcal{K} = \left\{ u \in L^2(Q_T); \|u\|_{L^\infty(0, T; L^2(\Omega))}^2 + \eta \|u\|_{L^2(0, T; H_0^1(\Omega))}^2 \leq A, \|\partial_t u\|_{L^2(0, T; H^{-1}(\Omega))} \leq B \right\},$$

The constants A and B will be fixed later. The set \mathcal{K} is a compact convex of $L^2(0, T; L^2(\Omega))$ (The compactness is due to the Aubin–Simon theorem [73]).

Let \mathcal{T} be a map from $L^2(0, T; L^2(\Omega))$ to $L^2(0, T; L^2(\Omega))$ defined by

$$\mathcal{T}(\bar{u}) = u,$$

where u is the unique solution to the following **linear** parabolic equation

$$\partial_t u - \operatorname{div}(a_\eta(\bar{u}) \nabla u - f(\bar{u}) \mathbf{V}) - g(\bar{u}) \operatorname{div}(\mathbf{V}) + a_\eta(\bar{u}) \nabla u \cdot \mathbf{V} = 0 \quad (5.19)$$

with the associate initial and boundary conditions. The existence of a unique solution to problem (5.19) is obtained using the Galerkin method [56, 28]. Indeed, there exists a unique solution u to problem (5.19) verifying : $u \in L^2(0, T; H_0^1(\Omega)) \cap \mathcal{C}^0(0, T; L^2(\Omega))$, $\partial_t u \in L^2(0, T; H^{-1}(\Omega))$ such that, we have the following weak formulation : $\forall \varphi \in L^2(0, T; H_0^1(\Omega))$,

$$\begin{aligned} \int_0^T \langle \partial_t u, \varphi \rangle dt + \int_{Q_T} a_\eta(\bar{u}) \nabla u \cdot \nabla \varphi dx dt - \int_{Q_T} f(\bar{u}) \mathbf{V} \cdot \nabla \varphi dx dt \\ + \int_{Q_T} g(\bar{u}) \mathbf{V} \cdot \nabla \varphi dx dt + \int_{Q_T} a_\eta(\bar{u}) \nabla u \cdot \mathbf{V} \varphi dx dt = 0. \end{aligned} \quad (5.20)$$

Lemma 5.8. \mathcal{T} is an application from \mathcal{K} to \mathcal{K} .

Proof. Since $u \in L^2(0, T; H_0^1(\Omega))$, one takes the solution u as a test function in the weak formulation (5.20), and gets, for all $t \in (0, T)$, that

$$E_1 + E_2 = E_3 + E_4, \quad (5.21)$$

where

$$\begin{aligned} E_1 &= \frac{1}{2} \int_0^t \left(\frac{d}{dt} \int_\Omega |u(\mathbf{x}, \tau)|^2 dx \right) d\tau = \frac{1}{2} \|u(t)\|_{L^2(\Omega)}^2 - \frac{1}{2} \|u_0\|_{L^2(\Omega)}^2, \\ E_2 &= \int_0^t \int_\Omega a_\eta(\bar{u}) \nabla u \cdot \nabla u dx d\tau, \\ E_3 &= \int_0^t \int_\Omega (f(\bar{u}) - g(\bar{u})) \nabla u \cdot \mathbf{V} dx d\tau, \\ E_4 &= - \int_0^t \int_\Omega a_\eta(\bar{u}) \nabla u \cdot \mathbf{V} u dx d\tau. \end{aligned}$$

From the degeneracy of the dissipation function a as well as its continuous extension, one has

$$E_2 = \int_0^t \int_\Omega (a(\bar{u}) + \eta) \nabla u \cdot \nabla u dx d\tau \geq \eta \int_0^t \int_\Omega \nabla u \cdot \nabla u dx d\tau. \quad (5.22)$$

Now, relying on the continuous extension of the functions f and g , the Cauchy-Schwarz, and the weighted Young inequality, one gets

$$\begin{aligned} |E_3| &\leq \int_0^t \int_\Omega |f(\bar{u}) - g(\bar{u})| |\nabla u \cdot \mathbf{V}| dx d\tau \leq C_{f,g} \int_0^t \int_\Omega |\nabla u \cdot \mathbf{V}| dx d\tau \\ &\leq C_{f,g} \|\nabla u\|_{(L^2(Q_t))^d} \cdot \|\mathbf{V}\|_{(L^2(Q_t))^d} \leq \kappa \|\nabla u\|_{(L^2(Q_t))^d}^2 + \frac{C_{f,g,Q_t}}{4\kappa} \|\mathbf{V}\|_{(L^\infty(Q_t))^d}^2, \end{aligned} \quad (5.23)$$

where κ is a constant to be specified later.

In the same manner, we have the following estimate

$$\begin{aligned} |E_4| &\leq \int_0^t \int_\Omega |a_\eta(\bar{u})| |\nabla u \cdot \mathbf{V} u| dx d\tau \leq C_{a,\eta} \int_0^t \int_\Omega |\nabla u \cdot \mathbf{V} u| dx d\tau \\ &\leq C_{a,\eta} \|\nabla u\|_{(L^2(Q_t))^d} \cdot \|\mathbf{V} u\|_{(L^2(Q_t))^d} \leq \kappa \|\nabla u\|_{(L^2(Q_t))^d}^2 + \frac{C_{f,g,\mathbf{V}}}{4\kappa} \|u\|_{(L^2(Q_t))^d}^2. \end{aligned} \quad (5.24)$$

Choosing the constant $\kappa = \frac{\eta}{4}$ and plugging estimates (5.22)–(5.24) into equation (5.21) one can conclude that

$$\|u(t)\|_{L^2(\Omega)}^2 + \eta \int_0^t \int_\Omega \nabla u \cdot \nabla u dx d\tau \leq \kappa_1 + \kappa_2 \|u\|_{L^2(Q_t)}^2, \quad (5.25)$$

or this is the same as

$$\|u(t)\|_{L^2(\Omega)}^2 + \eta \|\nabla u\|_{(L^2(Q_t))^d}^2 \leq \kappa_1 + \kappa_2 \int_0^t \|u(\tau)\|_{L^2(\Omega)}^2 dt, \quad (5.26)$$

where $\kappa_1 = \|u_0\|_{L^2(\Omega)}^2 + \frac{C_{f,g,Q_t}}{2\kappa} \|\mathbf{V}\|_{(L^\infty(Q_t))^d}^2$ and $\kappa_2 = \frac{C_{f,g,\mathbf{V}}}{4\kappa}$.

From estimate (5.26), and thanks to the lemma of Grönwall, one can deduce that there exists a constant $\kappa_3 = (\kappa_1 \exp(\kappa_2 t)) > 0$ such that

$$\|u\|_{L^2(Q_t)}^2 \leq \kappa_3, \quad \forall t \in (0, T). \quad (5.27)$$

Plugging estimate (5.27) into estimate (5.25), one has

$$\|u(t)\|_{L^2(\Omega)}^2 + \eta \|\nabla u\|_{(L^2(Q_t))^d}^2 \leq A, \quad \forall t \in (0, T),$$

where $A = \kappa_1 + \kappa_2 \kappa_3$.

Consequently, one deduces that

$$\|u\|_{L^\infty(0,T;L^2(\Omega))}^2 + \eta \|u\|_{L^2(0,T;H_0^1(\Omega))}^2 \leq A.$$

It remains to show the estimate on $\partial_t u$. To do this, we take $\varphi \in L^2(0, T; H_0^1(\Omega))$ as a test function into the weak formulation (5.20), one gets

$$\begin{aligned} \left| \int_0^T \langle \partial_t u, \varphi \rangle dt \right| &\leq \int_{Q_T} |f(\bar{u}) - g(\bar{u})| |\mathbf{V} \cdot \nabla \varphi| d\mathbf{x} dt + \int_{Q_T} |a_\eta \bar{u}| |\nabla u \cdot (\nabla \varphi + \mathbf{V} \varphi)| d\mathbf{x} dt \\ &\leq C_{f,g} \|\mathbf{V}\|_{(L^2(Q_T))^d} \|\nabla \varphi\|_{(L^2(Q_T))^d} + C_{a,\eta} \|\nabla u\|_{(L^2(Q_T))^d} \|\nabla \varphi\|_{(L^2(Q_T))^d} \\ &\quad + C_{a,\eta,\mathbf{V}} \|\nabla u\|_{(L^2(Q_T))^d} \|\varphi\|_{L^2(Q_T)}. \end{aligned}$$

Note that the Poincaré inequality implies the existence of a constant $\kappa_4 > 0$ (depending only on the domain Ω) such that

$$\|\varphi\|_{L^2(Q_T)} \leq \kappa_4 \|\nabla \varphi\|_{(L^2(Q_T))^d}.$$

Therefore, one can deduce that

$$\left| \int_0^t \langle \partial_t u, \varphi \rangle dt \right| \leq B \|\nabla \varphi\|_{(L^2(Q_T))^d}.$$

This ends the proof of the lemma. \square

Lemma 5.9. \mathcal{T} is a continuous application.

Proof. Let $(\bar{u}_n)_n$ be a sequence of \mathcal{K} and $\bar{u} \in \mathcal{K}$ such that $\bar{u}_n \rightarrow \bar{u}$ converges strongly in $L^2(0, T; L^2(\Omega))$. In order to prove the lemma, it suffices to show that

$$\mathcal{T}(\bar{u}_n) = u_n \rightarrow \mathcal{T}(\bar{u}) = u \text{ converges strongly in } L^2(0, T; L^2(\Omega)).$$

For all $\varphi \in L^2(0, T; H_0^1(\Omega))$, the sequence $(u_n)_n$ satisfies

$$\begin{aligned} \int_0^T \langle \partial_t u_n, \varphi \rangle dt + \int_{Q_T} a_\eta(\bar{u}_n) \nabla u_n \cdot \nabla \varphi d\mathbf{x} dt - \int_{Q_T} f(\bar{u}_n) \mathbf{V} \cdot \nabla \varphi d\mathbf{x} dt \\ + \int_{Q_T} g(\bar{u}_n) \mathbf{V} \cdot \nabla \varphi d\mathbf{x} dt + \int_{Q_T} a_\eta(\bar{u}_n) \mathbf{V} \cdot \nabla u_n \varphi d\mathbf{x} dt = 0. \end{aligned} \quad (5.28)$$

Let us denote v_n by $v_n = u_n - u$ and substrat equation (5.20) from equation (5.28), one has $\forall \varphi \in L^2(0, T; H_0^1(\Omega))$

$$\begin{aligned} \int_0^T \langle \partial_t v_n, \varphi \rangle dt + \int_{Q_T} a_\eta(\bar{u}_n) \nabla v_n \cdot \nabla \varphi d\mathbf{x} dt + \int_{Q_T} (a_\eta(\bar{u}_n) - a_\eta(\bar{u})) \nabla u \cdot \nabla \varphi d\mathbf{x} dt \\ - \int_{Q_T} (f(\bar{u}_n) - f(\bar{u})) \mathbf{V} \cdot \nabla \varphi d\mathbf{x} dt + \int_{Q_T} (g(\bar{u}_n) - g(\bar{u})) \mathbf{V} \cdot \nabla \varphi d\mathbf{x} dt \\ + \int_{Q_T} a_\eta(\bar{u}_n) \mathbf{V} \cdot \nabla v_n \varphi d\mathbf{x} dt + \int_{Q_T} (a_\eta(\bar{u}_n) - a_\eta(\bar{u})) \nabla u \cdot \mathbf{V} \varphi d\mathbf{x} dt = 0. \end{aligned} \quad (5.29)$$

Now, we take $\varphi = v_n$ as a test function in equation (5.29), and a parameter $\delta > 0$ defined later, we have the following equation

$$\sum_{1 \leq i \leq 7} H_i = 0,$$

where

$$\begin{aligned} H_1 &= \int_0^t \langle \partial_t v_n, v_n \rangle d\tau = \frac{1}{2} \|v_n(t)\|_{L^2(\Omega)}^2, \\ H_2 &= \int_{Q_t} a_\eta(\bar{u}_n) \nabla v_n \cdot \nabla v_n d\tau d\mathbf{x} \geq \eta \|\nabla v_n\|_{(L^2(Q_t))^d}^2, \\ H_3 &= \left| \int_{Q_t} (a_\eta(\bar{u}_n) - a_\eta(\bar{u})) \nabla u \cdot \nabla v_n d\tau d\mathbf{x} \right| \\ &\leq \delta \|\nabla v_n\|_{(L^2(Q_t))^d}^2 + \frac{1}{4\delta} \|(a_\eta(\bar{u}_n) - a_\eta(\bar{u})) \nabla u\|_{(L^2(Q_t))^d}^2, \\ H_4 &= \left| \int_{Q_t} (f(\bar{u}_n) - f(\bar{u})) \nabla v_n \cdot \mathbf{V} d\tau d\mathbf{x} \right| \\ &\leq \delta \|\nabla v_n\|_{(L^2(Q_t))^d}^2 + \frac{1}{4\delta} \|(f(\bar{u}_n) - f(\bar{u})) \mathbf{V}\|_{(L^2(Q_t))^d}^2, \\ H_5 &= \left| \int_{Q_t} (g(\bar{u}_n) - g(u_n)) \nabla v_n \cdot \mathbf{V} d\tau d\mathbf{x} \right| \\ &\leq \delta \|\nabla v_n\|_{(L^2(Q_t))^d}^2 + \frac{1}{4\delta} \|(g(\bar{u}_n) - g(\bar{u})) \mathbf{V}\|_{(L^2(Q_t))^d}^2, \\ H_6 &= \left| \int_{Q_t} a_\eta(\bar{u}_n) \nabla v_n \cdot \mathbf{V} v_n d\tau d\mathbf{x} \right| \leq \delta \|\nabla v_n\|_{(L^2(Q_t))^d}^2 + \frac{C_{a,\eta}}{4\delta} \|v_n\|_{L^2(Q_t)}^2, \\ H_7 &= \left| \int_{Q_t} (a_\eta(\bar{u}_n) - a_\eta(\bar{u})) \nabla u \cdot \mathbf{V} v_n d\tau d\mathbf{x} \right| \\ &\leq \delta \|\nabla v_n\|_{(L^2(Q_t))^d}^2 + C_{a,\eta,\mathbf{V}} \|(a_\eta(\bar{u}_n) - a_\eta(\bar{u})) \nabla u\|_{(L^2(Q_t))^d}^2. \end{aligned}$$

Plugging these estimates into equation (5.29) and choosing $\delta = \frac{\eta}{12}$, one can deduce that

$$\|v_n(t)\|_{L^2(\Omega)}^2 \leq \kappa_5 + \frac{6C_{a,\eta}}{\eta} \|v_n\|_{L^2(Q_t)}^2,$$

where

$$\begin{aligned} \kappa_5 &= \frac{1}{2\delta} \left(\|(a_\eta(\bar{u}_n) - a_\eta(\bar{u})) \nabla u\|_{(L^2(Q_t))^d}^2 + \|(f(\bar{u}_n) - f(\bar{u})) \mathbf{V}\|_{(L^2(Q_t))^d}^2 \right. \\ &\quad \left. \|(g(\bar{u}_n) - g(\bar{u})) \mathbf{V}\|_{(L^2(Q_t))^d}^2 + C_{a,\eta,\mathbf{V}} \|(a_\eta(\bar{u}_n) - a_\eta(\bar{u})) \nabla u\|_{(L^2(Q_t))^d}^2 \right). \end{aligned}$$

Now, thanks to the Grönwall lemma, one can deduce that

$$\|v_n(t)\|_{L^2(\Omega)}^2 \leq \kappa_5 \exp\left(\frac{6C_{a,\eta}t}{\eta}\right).$$

Note that κ_5 tends to zero as $n \rightarrow \infty$. Indeed, one can get the result using Lebesgue's dominated convergence theorem, the uniqueness of a solution to problem (5.19), and using the continuity of each of the functions a , f , and g . Therefore, one gets

$$\|v_n(t)\|_{L^2(\Omega)} \xrightarrow{n \rightarrow +\infty} 0, \quad \forall t \in (0, T).$$

In other words, one has

$$\|v_n\|_{L^\infty(0,T;L^2(\Omega))} \rightarrow 0 \text{ quand } n \rightarrow +\infty,$$

which implies that

$$u_n \rightarrow u \text{ strongly in } L^2(0, T; L^2(\Omega)).$$

This ends the proof of this lemma. \square

Using previous results, Green-Riemann's theorem, and Schauder's fixed-point theorem, one can deduce that there exists at least one solution to the nondegenerate problem (5.14)–(5.16) in the sense of theorem 5.7. It remains to show that the solution verifies the maximum principle.

5.3.2 Maximum principle on the saturation

In this section, we aim to prove that the solution to the nondegenerate problem (5.14)–(5.16) is stable in the sense of verifying the maximum principle. Specifically, we have the following lemma.

Lemma 5.10. *Let u be a solution to the nondegenerate system (5.14)–(5.16) under the assumptions (H1) – (H3) and (H5). Then, the solution u satisfies*

$$0 \leq u(\mathbf{x}, t) \leq 1, \quad \text{for a.e. } (\mathbf{x}, t) \in Q_T.$$

Proof. Let u^- be the function defined by $u^- = \max(-u, 0) = \frac{|u| - u}{2} \geq 0$. Stampacchia's Theorem ensures that $u^- \in L^2(0, T; H_0^1(\Omega))$ since $u \in L^2(0, T; H_0^1(\Omega))$. Therefore, one can consider $-u^-$ as a test function into the weak formulation (5.20), and gets using the inequality $u = u^+ - u^-$ that

$$\begin{aligned} \int_0^t \langle \partial_t u^-, u^- \rangle dt + \int_{Q_t} a_\eta(u) \nabla u^- \cdot \nabla u^- d\mathbf{x} dt + \int_{Q_t} f(u) \mathbf{V} \cdot \nabla u^- d\mathbf{x} dt \\ + \int_{Q_t} g(u) \operatorname{div}(\mathbf{V}) u^- d\mathbf{x} dt + \int_{Q_t} a_\eta(u) \nabla u^- \cdot \mathbf{V} u^- d\mathbf{x} dt = 0. \end{aligned} \quad (5.30)$$

For the first term of equation (5.30), one has

$$\int_0^t \langle \partial_t u^-, u^- \rangle dt = \frac{1}{2} \int_0^t \frac{d}{dt} \|u^-(\tau)\|_{L^2(\Omega)}^2 d\tau = \frac{1}{2} \|u^-(t)\|_{L^2(\Omega)}^2,$$

since $u_0^- = 0$ due to the nonnegativity of the function $u_0(\mathbf{x})$ for a.e. $\mathbf{x} \in \Omega$.

Now, we use the definition of the function a_η and the degeneracy of the dissipation a to conclude that

$$\int_{Q_t} a_\eta(u) \nabla u^- \cdot \nabla u^- d\mathbf{x} dt \geq \eta \int_{Q_t} \nabla u^- \cdot \nabla u^- d\mathbf{x} dt = \eta \|\nabla u^-\|_{(L^2(Q_t))^d}^2. \quad (5.31)$$

Furthermore, we rely on the continuous extension by zero of the functions $f(u)$ and $g(u)$ for $u \leq 0$, to deduce that the third and the fourth terms in equation (5.30) are equal to zero. Let us now focus on the last term of equation (5.30). Indeed, by the Cauchy-Schwarz inequality as well as the weighted Young inequality, one has

$$\int_{Q_t} a_\eta(u) \nabla u^- \cdot \mathbf{V} u^- \, d\mathbf{x} \, dt \leq \frac{\eta}{2} \|\nabla u^-\|_{(L^2(Q_t))^d}^2 + \frac{C_{a,\eta,\mathbf{V}}}{2} \int_0^t \|u^-(\tau)\|_{L^2(\Omega)}^2 \, d\tau. \quad (5.32)$$

Substituting estimates (5.31)–(5.32) into equation (5.30), this yields

$$\|u^-(t)\|_{L^2(\Omega)}^2 + \eta \|\nabla u^-\|_{(L^2(Q_t))^d}^2 \leq C_{a,\eta,\mathbf{V}} \int_0^t \|u^-(\tau)\|_{L^2(\Omega)}^2 \, d\tau.$$

Denote by ϕ the function defined by $\phi(t) = \int_0^t \|u^-(\tau)\|_{L^2(\Omega)}^2 \, d\tau$, then one has

$$\frac{d\phi}{dt}(t) \leq C_{a,\eta,\mathbf{V}} \phi(t).$$

Applying, the Grönwall lemma, one can deduce that

$$\|u^-(t)\|_{L^2(\Omega)}^2 = \phi(t) \leq \phi(0) \exp(C_{a,\eta,\mathbf{V}} t) = 0.$$

As a consequence, $u^-(\mathbf{x}, t) = 0$, for a.e. $(\mathbf{x}, t) \in Q_T$, i.e. $u(\mathbf{x}, t) \geq 0$, for a.e. $(\mathbf{x}, t) \in Q_T$.

It remains to show that $u(\mathbf{x}, t) \leq 1$, for a.e. $(\mathbf{x}, t) \in Q_T$. To do this, it suffices to prove that $(u - 1)^+ = 0$. Thus, we multiply the saturation equation (5.14) by the regular function $(u - 1)^+ \in L^2(0, T; H_0^1(\Omega))$ and integrate the resulting equation over $\Omega \times (0, t)$, this yields

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_0^t \|(u - 1)^+(\tau)\|_{L^2(\Omega)}^2 \, d\tau + \int_{Q_t} a_\eta(u) \nabla (u - 1)^+ \cdot \nabla (u - 1)^+ \, d\mathbf{x} \, dt \\ & - \int_{Q_t} f(u) \mathbf{V} \cdot \nabla (u - 1)^+ \, d\mathbf{x} \, dt - \int_{Q_t} g(u) \operatorname{div}(\mathbf{V}) (u - 1)^+ \, d\mathbf{x} \, dt \\ & + \int_{Q_t} a_\eta(u) \nabla u \cdot \mathbf{V} (u - 1)^+ \, d\mathbf{x} \, dt = 0. \end{aligned} \quad (5.33)$$

Now, we proceed as before and get the estimates for each term of equation (5.33).

For the first term, one has

$$\frac{1}{2} \frac{d}{dt} \int_0^t \|(u - 1)^+(\tau)\|_{L^2(\Omega)}^2 \, d\tau = \frac{1}{2} \|(u - 1)^+(t)\|_{L^2(\Omega)}^2 - \frac{1}{2} \|(u_0 - 1)^+\|_{L^2(\Omega)}^2, \quad (5.34)$$

and since $u_0 \leq 1$ then the second term on the right-hand side of (5.34) is equal to zero.

For the third and the fourth term of equation (5.33), by using the fact that $f(u) = g(u) = 1$ for all $u \geq 1$ and the fact that $\mathbf{V} \cdot \mathbf{n} \leq 0$ on $\partial\Omega$, one has

$$\begin{aligned} & - \int_{Q_t} f(u) \mathbf{V} \cdot \nabla (u - 1)^+ \, d\mathbf{x} \, dt - \int_{Q_t} g(u) \operatorname{div}(\mathbf{V}) (u - 1)^+ \, d\mathbf{x} \, dt \\ & = - \int_{\Sigma_T} (u - 1)^+ \mathbf{V} \cdot \mathbf{n} \, d\sigma(\mathbf{x}) \, dt \geq 0 \end{aligned}$$

Finally, for the last term of equation (5.34), we use again the extension by $a(1)$ of the dissipation function a for $u > 1$, the Cauchy-Schwarz inequality and the weighted Young inequality, and get

the following estimate

$$\begin{aligned} \int_{Q_t} a_\eta(u) \nabla u \cdot \mathbf{V}(u-1)^+ d\mathbf{x} dt &= \int_{Q_t} a_\eta(u) \nabla(u-1) \cdot \mathbf{V}(u-1)^+ d\mathbf{x} dt \\ &\leq \frac{\eta}{2} \|\nabla(u-1)^+\|_{(L^2(Q_t))^2}^2 + \frac{C_{a,\eta,\mathbf{V}}}{2} \int_0^t \|(u-1)^+(\tau)\|_{L^2(\Omega)}^2 d\tau. \end{aligned}$$

Plugging the previous estimates into equation (5.33), one has

$$\|(u-1)^+(t)\|_{L^2(\Omega)}^2 + \eta \|\nabla(u-1)^+\|_{(L^2(Q_t))^d}^2 \leq C_{a,\eta,\mathbf{V}} \int_0^t \|(u-1)^+(\tau)\|_{L^2(\Omega)}^2 d\tau.$$

One can conclude, using the Grönwall lemma, that $u(\mathbf{x}, t) \leq 1$, for a.e. $(\mathbf{x}, t) \in Q_T$. This ends the proof of lemma 5.10. \square

The proof of theorem 5.7 is now completed. \square

5.4 Proof of theorem 5.2.

In the previous section, we have shown that the nondegenerate system (5.14)–(5.16) admits at least one weak solution. Here, we are going to prove theorem (5.2), the proof is based on the establishment of estimates on the solutions independent of the parameter η , and next on the passage to the limit as η tends to zero.

From the definition (5.4) of k , we have

$$k(u) = k(u_*) \exp \left(\int_{u_*}^u (f(\tau) - g(\tau))^{-1} g'(\tau) d\tau \right), \quad \text{for all } u \geq u_*.$$

As a consequence of assumption (H4a), there exist two constants c_3 and c_4 depending on f , g , and u_* such that

$$c_3(1-u)^{-1} \leq k(u) \leq c_4(1-u)^{-1}, \quad \forall u_* \leq u < 1. \quad (5.35)$$

Indeed, we have

$$k(u_*) \exp \left(c_1 C_{g'} \int_{u_*}^{u_\eta} \frac{1}{1-\tau} d\tau \right) \leq k(u_\eta) \leq k(u_*) \exp \left(c_2 \|g'\|_\infty \int_{u_*}^{u_\eta} \frac{1}{1-\tau} d\tau \right).$$

That is

$$\frac{c_1 C_{g'} k(u_*) (1-u_*)}{1-u_\eta} \leq k(u_\eta) \leq \frac{c_2 \|g'\|_\infty k(u_*) (1-u_*)}{1-u_\eta}.$$

Denoting by $c_3 = c_1 C_{g'} k(u_*) (1-u_*)$ and $c_4 = c_2 \|g'\|_\infty k(u_*) (1-u_*)$, then one obtains the confinement (5.35).

Now, using the confinement of the function k and denoting by $c_5 = c_1 c_3 C_{g'}$ and $c_6 = c_2 c_4 \|g'\|_\infty$, one can easily obtain that

$$\frac{c_5}{(1-u_\eta)^{-2}} \leq k'(u_\eta) \leq \frac{c_6}{(1-u_\eta)^{-2}}.$$

Lemma 5.11. *Let $L(u_0)$ belongs to $L^1(\Omega)$. Under the assumptions (H1) – (H3) and (H4a), the weak solution to the nondegenerate system (5.14)–(5.16) verifies*

$$\text{The sequence } (\nabla u_\eta)_\eta \text{ is uniformly bounded in } (L^2(Q_T))^d. \quad (5.36)$$

$$\text{The sequence } \left(\sqrt{\eta k'(u_\eta)} \nabla u_\eta \right)_\eta \text{ is uniformly bounded in } (L^2(Q_T))^d. \quad (5.37)$$

$$\text{The sequence } (\partial_t u_\eta)_\eta \text{ is uniformly bounded in } L^2(0, T; H^{-1}(\Omega)). \quad (5.38)$$

$$\text{The sequence } (u_\eta)_\eta \text{ relatively compact in } L^2(Q_T). \quad (5.39)$$

Proof. To prove (5.36) and (5.37), we multiply the saturation equation (5.14) by k and integrate it over Ω , one gets

$$\begin{aligned} & \frac{d}{dt} \int_{\Omega} L(u_{\eta}) \, d\mathbf{x} + \int_{\Omega} a(u_{\eta}) k'(u_{\eta}) \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x} + \eta \int_{\Omega} k'(u_{\eta}) \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x} \\ &= \int_{\Omega} (f(u_{\eta}) - g(u_{\eta})) \mathbf{V} \cdot \nabla k(u_{\eta}) \, d\mathbf{x} - \int_{\Omega} g'(u_{\eta}) k(u_{\eta}) \mathbf{V} \cdot \nabla u_{\eta} \, d\mathbf{x} \quad (5.40) \\ & \quad - \int_{\Omega} a(u_{\eta}) k(u_{\eta}) \mathbf{V} \cdot \nabla u_{\eta} \, d\mathbf{x} - \eta \int_{\Omega} k(u_{\eta}) \mathbf{V} \cdot \nabla u_{\eta} \, d\mathbf{x}. \end{aligned}$$

We denote $\Omega_1 = \Omega \cap \{u_{\eta} < u_*\}$ and $\Omega_2 = \Omega \cap \{u_{\eta} \geq u_*\}$; then the whole integral on Ω can be split into two sub integrals on Ω_1 and Ω_2 respectively. Within the region Ω_1 , we have $a(u_{\eta}) \geq a_0$ and $k'(u_{\eta}) = 1$, this yields

$$\int_{\Omega_1} a(u_{\eta}) k'(u_{\eta}) \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x} \geq a_0 \int_{\Omega_1} \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x} = a_0 \|\nabla u_{\eta}\|_{(L^2(\Omega))^d}^2.$$

Furthermore, we have the following estimates

$$\begin{aligned} & \left| \int_{\Omega_1} (f(u_{\eta}) - g(u_{\eta})) k'(u_{\eta}) \mathbf{V} \cdot \nabla u_{\eta} \, d\mathbf{x} \right| \leq \int_{\Omega_1} |(f(u_{\eta}) - g(u_{\eta})) \mathbf{V} \cdot \nabla u_{\eta}| \, d\mathbf{x} \\ & \leq \frac{a_0}{6} \|\nabla u_{\eta}\|_{(L^2(\Omega_1))^d}^2 + \frac{3C}{2a_0} \|\mathbf{V}\|_{(L^2(\Omega_1))^d}^2, \\ & \left| \int_{\Omega_1} (g'(u_{\eta}) + a(u_{\eta})) k(u_{\eta}) \mathbf{V} \cdot \nabla u_{\eta} \, d\mathbf{x} \right| \leq \int_{\Omega_1} |(g'(u_{\eta}) + a(u_{\eta})) \mathbf{V} \cdot \nabla u_{\eta}| \, d\mathbf{x} \\ & \leq \frac{a_0}{3} \|\nabla u_{\eta}\|_{(L^2(\Omega_1))^d}^2 + \frac{3(\|g'\|_{\infty} + \|a\|_{\infty})^2}{4a_0} \|\mathbf{V}\|_{(L^2(\Omega_1))^d}^2, \\ & \left| \eta \int_{\Omega_1} k(u_{\eta}) \mathbf{V} \cdot \nabla u_{\eta} \, d\mathbf{x} \right| \leq \int_{\Omega_1} |\eta u_{\eta} \mathbf{V} \cdot \nabla u_{\eta}| \, d\mathbf{x} \leq \int_{\Omega_1} |\eta \mathbf{V} \cdot \nabla u_{\eta}| \, d\mathbf{x} \\ & \leq \frac{1}{2} \|\mathbf{V}\|_{(L^2(\Omega_1))^d}^2 + \frac{1}{2} \int_{\Omega_1} \eta k'(u_{\eta}) \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x}. \end{aligned}$$

In region Ω_2 , we have

$$\begin{aligned} & \int_{\Omega_2} ((f(u_{\eta}) - g(u_{\eta})) k'(u_{\eta}) - g'(u_{\eta}) k(u_{\eta})) \mathbf{V} \cdot \nabla u_{\eta} \, d\mathbf{x} = 0, \\ & \int_{\Omega_2} a(u_{\eta}) k'(u_{\eta}) \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x} \geq c_5 \int_{\Omega_2} a(u_{\eta}) (1 - u_{\eta})^{-2} \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x}. \end{aligned}$$

The following estimates hold

$$\begin{aligned} & \left| \int_{\Omega_2} \eta k(u_{\eta}) \mathbf{V} \cdot \nabla u_{\eta} \, d\mathbf{x} \right| \leq c_4 \int_{\Omega_2} |\sqrt{\eta} (1 - u_{\eta})^{-1} \mathbf{V} \cdot \nabla u_{\eta}| \, d\mathbf{x} \\ & \leq \frac{c_4^2}{2c_5} \|\mathbf{V}\|_{(L^2(\Omega))^d}^2 + \frac{1}{2} \int_{\Omega_2} \eta k'(u_{\eta}) \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x}, \\ & \left| \int_{\Omega_2} a(u_{\eta}) k(u_{\eta}) \mathbf{V} \cdot \nabla u_{\eta} \, d\mathbf{x} \right| \leq c_4 \int_{\Omega_2} \left| \sqrt{a(u_{\eta})} (1 - u_{\eta})^{-1} \mathbf{V} \cdot \sqrt{a(u_{\eta})} \nabla u_{\eta} \right| \, d\mathbf{x} \\ & \leq c_4 \left(\int_{\Omega_2} a(u_{\eta}) \mathbf{V}^2 \, d\mathbf{x} \right)^{\frac{1}{2}} \left(\int_{\Omega_2} a(u_{\eta}) (1 - u_{\eta})^{-2} \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x} \right)^{\frac{1}{2}} \\ & \leq \frac{c_4^2 \|a\|_{\infty}^2}{2c_5} \|\mathbf{V}\|_{(L^2(\Omega))^d}^2 + \frac{c_5}{2} \int_{\Omega_2} a(u_{\eta}) (1 - u_{\eta})^{-2} \nabla u_{\eta} \cdot \nabla u_{\eta} \, d\mathbf{x}. \end{aligned}$$

From the degeneracy assumption (H4a), we have $a(u_\eta) \geq m_1(1 - u_\eta)^{r_2}$ for an exponent r_2 such that $r_2 \leq 2$, and therefore

$$\int_{\Omega_2} a(u_\eta) (1 - u_\eta)^{-2} \nabla u_\eta \cdot \nabla u_\eta \, d\mathbf{x} \geq m_1(1 - u_*)^{r_2-2} \int_{\Omega_2} \nabla u_\eta \cdot \nabla u_\eta \, d\mathbf{x}.$$

Plugging the previous estimates into equation (5.40), this yields

$$\begin{aligned} \frac{d}{dt} \int_{\Omega} L(u_\eta) \, d\mathbf{x} + \frac{a_0}{2} \int_{\Omega_1} \nabla u_\eta \cdot \nabla u_\eta \, d\mathbf{x} + \frac{c_5 m_1 (1 - u_*)^{r_2-2}}{2} \int_{\Omega_2} \nabla u_\eta \cdot \nabla u_\eta \, d\mathbf{x} \\ + \frac{1}{2} \int_{\Omega_2} \eta k'(u_\eta) \nabla u_\eta \cdot \nabla u_\eta \, d\mathbf{x} \leq C_{a,g',a_0}, \end{aligned}$$

where, C_{a,g',a_0} is a nonnegative constant equals to

$$C_{a,g',a_0} = \frac{1}{2} \left(\frac{3C}{a_0} + \frac{3(\|g'\|_\infty + \|a\|_\infty)}{2a_0} + \frac{c_4^2 \|a\|_\infty^2}{c_5} + 1 + \frac{c_4^2}{c_5} \right) \|\mathbf{V}\|_{(L^2(\Omega))^d}^2.$$

As a consequence, one has

$$2 \frac{d}{dt} \int_{\Omega} L(u_\eta) \, d\mathbf{x} + C_{a_0,u_*,m_1} \|\nabla u_\eta\|_{L^2(\Omega)}^2 + \left\| \sqrt{\eta k'(u_\eta)} \nabla u_\eta \right\|_{L^2(\Omega)}^2 \leq 2C_{a,g',a_0}, \quad (5.41)$$

where, $C_{a_0,u_*,m_1} = \min(a_0, c_5 m_1 (1 - u_*)^{r_2-2})$.

Non, we integrate inequality (5.41) with respect to the time over $(0, t)$, $t \in (0, T)$, one gets

$$\begin{aligned} 2 \int_{\Omega} L(u_\eta) \, d\mathbf{x} + C_{a_0,u_*} \|\nabla u_\eta\|_{L^2(Q_t)}^2 + \left\| \sqrt{\eta k'(u_\eta)} \nabla u_\eta \right\|_{L^2(Q_t)}^2 \\ \leq 2 \left(TC_{a,g',a_0} + \|L(u_0)\|_{L^1(\Omega)} \right). \end{aligned}$$

One can conclude, using these estimate, that the sequences $(\nabla u_\eta)_\eta$ and $(\sqrt{\eta k'(u_\eta)} \nabla u_\eta)_\eta$ are uniformly bounded in $(L^2(0, T; L^2(\Omega)))^d$.

Let us prove the third part (5.38) of this lemma, for that we take $\varphi \in L^2(0, T; H_0^1(\Omega))$ as a test function into the weak formulation (5.18), one gets

$$\begin{aligned} \left| \int_0^T \langle \partial_t u_\eta, \varphi \rangle \, dt \right| \leq \int_{Q_T} |a(u_\eta) \nabla u_\eta \cdot \nabla \varphi| \, d\mathbf{x} \, dt + \int_{Q_T} |a(u_\eta) \mathbf{V} \cdot \nabla u_\eta \varphi| \, d\mathbf{x} \, dt \\ + \int_{Q_T} |(f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla \varphi| \, d\mathbf{x} \, dt + \int_{Q_T} |g'(u_\eta) \mathbf{V} \cdot \nabla u_\eta \varphi| \, d\mathbf{x} \, dt \\ + \int_{Q_T} |\eta \nabla u_\eta \cdot \nabla \varphi| \, d\mathbf{x} \, dt + \int_{Q_T} |\eta \mathbf{V} \cdot \nabla u_\eta \varphi| \, d\mathbf{x} \, dt. \end{aligned}$$

Using the Cauchy-Schwarz inequality as well as the previous estimates, one can deduce that

$$\begin{aligned} \int_{Q_T} |a(u_\eta) \nabla u_\eta \cdot \nabla \varphi| \, d\mathbf{x} \, dt + \int_{Q_T} |(f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla \varphi| \, d\mathbf{x} \, dt \\ + \int_{Q_T} |\eta \nabla u_\eta \cdot \nabla \varphi| \, d\mathbf{x} \, dt \leq C_{a_0,a,u_*,g',L,\mathbf{V}} \|\nabla \varphi\|_{(L^2(Q_T))^d}. \end{aligned}$$

also, one gets

$$\begin{aligned} \int_{Q_T} |g'(u_\eta) \mathbf{V} \cdot \nabla u_\eta \varphi| \, d\mathbf{x} \, dt + \int_{Q_T} |a(u_\eta) \mathbf{V} \cdot \nabla u_\eta \varphi| \, d\mathbf{x} \, dt \\ + \int_{Q_T} |\eta \mathbf{V} \cdot \nabla u_\eta \varphi| \, d\mathbf{x} \, dt \leq C_{a_0, a, u_*, g', L, \mathbf{V}} \|\varphi\|_{L^2(Q_T)}. \end{aligned}$$

Using the Poincaré inequality, one deduces that there exists a constant $C_{a_0, a, u_*, g', L, \mathbf{V}, \Omega, T} > 0$ such that

$$|\langle \partial_t u_\eta, \varphi \rangle| \leq C_{a_0, a, u_*, g', L, \mathbf{V}, \Omega, T} \|\varphi\|_{L^2(0, T; H_0^1(\Omega))}, \quad \forall \varphi \in L^2(0, T; H_0^1(\Omega)).$$

This proves that the sequence $(\partial_t u_\eta)_\eta$ is uniformly bounded in $L^2(0, T; H^{-1}(\Omega))$.

It remains to show the last part (5.39) of the lemma. Indeed, using statements (5.36) and (5.38), we remark that the sequence $(u_\eta)_\eta$ belongs to the Sobolev space

$$\mathcal{W} = \{u_\eta; u_\eta \in L^2(0, T; H_0^1(\Omega)) \text{ and } \partial_t u_\eta \in L^2(0, T; H^{-1}(\Omega))\}.$$

Since $H_0^1(\Omega)$ is compactly embedded in $L^2(\Omega)$ and $L^2(\Omega)$ is continuously embedded in $H^{-1}(\Omega)$, then thanks to the Aubin–Lions lemma, \mathcal{W} is compactly embedded in $L^2(0, T; L^2(\Omega))$. Consequently, the sequence $(u_\eta)_\eta$ is relatively compact in $L^2(0, T; L^2(\Omega))$. This ends the proof of lemma 5.11. \square

The next step of the proof of theorem 5.2 is to pass to the limit as η tends to zero in the weak formulation (5.18).

Lemma 5.12. *We have the following convergences as η goes to zero*

$$u_\eta \longrightarrow u \text{ strongly in } L^2(Q_T), \quad (5.42)$$

$$\nabla u_\eta \longrightarrow \nabla u \text{ weakly in } (L^2(Q_T))^d, \quad (5.43)$$

Furthermore,

$$u_\eta \longrightarrow u \text{ almost everywhere in } Q_T, \quad (5.44)$$

$$0 \leq u(\mathbf{x}, t) \leq 1 \text{ almost everywhere in } Q_T, \quad (5.45)$$

$$\partial_t u_\eta \longrightarrow \partial_t u \text{ weakly in } L^2(0, T; H^{-1}(\Omega)). \quad (5.46)$$

Proof. To prove the strong convergence (5.42), we deduce from lemma 5.11 and thanks to the Aubin–Lions lemma that there exists a subsequence also denoted $(u_\eta)_\eta$ and such that

$$u_\eta \longrightarrow u \text{ strongly in } L^2(Q_T).$$

On the other hand, and thanks to the uniform bound of the sequence $(\nabla u_\eta)_\eta$ in $(L^2(Q_T))^d$, one can extract a subsequence such that

$$\nabla u_\eta \longrightarrow \nabla u \text{ weakly in } (L^2(Q_T))^d,$$

to which the identification of the limit is made using the classical properties on the distribution theory.

The convergence (5.44) is a consequence of the strong convergence (5.42), while the convergence (5.45) is obtained using the fact that $\{u \in L^2(Q_T), 0 \leq u(\mathbf{x}, t) \leq 1 \text{ for a.e. } (\mathbf{x}, t) \in Q_T\}$ is a closed subset of $L^2(Q_T)$.

Finally, the weak convergence (5.46) is a consequence of (5.38) and the strong convergence (5.42). \square

Using the convergences given in lemma 5.12 as well as the Lebesgue theorem, one has for all $\varphi \in L^2(0, T; H_0^1(\Omega))$

$$\int_{Q_T} \partial_t u_\eta \varphi \, d\mathbf{x} \, dt \longrightarrow \int_{Q_T} \partial_t u \varphi \, d\mathbf{x} \, dt, \text{ as } \eta \rightarrow 0. \quad (5.47)$$

$$\int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla u_\eta \, d\mathbf{x} \, dt \longrightarrow \int_{Q_T} a(u) \nabla u \cdot \nabla u \, d\mathbf{x} \, dt, \text{ as } \eta \rightarrow 0. \quad (5.48)$$

$$\int_{Q_T} (f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla \varphi \, d\mathbf{x} \, dt \longrightarrow \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla \varphi \, d\mathbf{x} \, dt, \text{ as } \eta \rightarrow 0. \quad (5.49)$$

Furthermore, one gets

$$\begin{aligned} \int_{Q_T} (g'(u_\eta) + a(u_\eta)) \mathbf{V} \cdot \nabla u_\eta \varphi \, d\mathbf{x} \, dt \\ \longrightarrow \int_{Q_T} (g'(u) + a(u)) \mathbf{V} \cdot \nabla u \varphi \, d\mathbf{x} \, dt \text{ as } \eta \rightarrow 0. \end{aligned} \quad (5.50)$$

It remains to show that

$$\int_{Q_T} \eta \nabla u_\eta \cdot \nabla \varphi \, d\mathbf{x} \, dt + \int_{Q_T} \eta \mathbf{V} \varphi \cdot \nabla u_\eta \, d\mathbf{x} \, dt \longrightarrow 0, \text{ as } \eta \rightarrow 0.$$

Indeed, using the Cauchy-Schwarz inequality and the uniform bound of $(\nabla u_\eta)_\eta$, one has

$$\int_{Q_T} \eta \nabla u_\eta \cdot \nabla \varphi \, d\mathbf{x} \, dt \leq \eta \|\nabla u_\eta\|_{(L^2(Q_T))^2} \|\nabla \varphi\|_{(L^2(Q_T))^2} \longrightarrow 0, \text{ as } \eta \rightarrow 0. \quad (5.51)$$

$$\int_{Q_T} \eta \mathbf{V} \varphi \cdot \nabla u_\eta \, d\mathbf{x} \, dt \leq \eta C_{\mathbf{V}} \|\nabla u_\eta\|_{(L^2(Q_T))^2} \|\nabla u_\eta\|_{L^2(Q_T)} \longrightarrow 0, \text{ as } \eta \rightarrow 0. \quad (5.52)$$

Using the convergence results (5.47)–(5.52), we can then pass to the limit in the weak formulation (5.18) and obtain the classical weak formulation on the limit solution u . This ends the proof of theorem 5.2.

5.5 Proof of theorem 5.4

In this section and under the degeneracy assumption (H4b) on the dissipation function a , we carry out some estimates on the solutions to the nondegenerate system (5.14)–(5.16) independent of η . We recall that the function β is defined as $\beta(u) = u^{r-1}$ and that h is its primitive defined by $h(u) = \int_0^u \beta(\mathbf{y}) \, d\mathbf{y}$.

Lemma 5.13. *Under the assumptions (H1) – (H3) and (H4b), the solutions to the nondegenerate system (5.14)–(5.16) satisfy*

$$0 \leq u_\eta(\mathbf{x}, t) \leq 1, \text{ for a.e. } (\mathbf{x}, t) \in Q_T. \quad (5.53)$$

$$\left(\sqrt{u_\eta^{r-2} a(u_\eta)} \nabla u_\eta \right)_\eta \text{ and } (u_\eta^{r-1} \nabla u_\eta)_\eta \text{ are uniformly bounded in } (L^2(Q_T))^d. \quad (5.54)$$

$$\left(\sqrt{\eta u_\eta^{r-2}} \nabla u_\eta \right)_\eta \text{ and } (a(u_\eta) \nabla u_\eta)_\eta \text{ are uniformly bounded in } (L^2(Q_T))^d. \quad (5.55)$$

$$(h(u_\eta))_\eta \text{ is uniformly bounded in } L^\infty(0, T; L^2(\Omega)). \quad (5.56)$$

$$(\partial_t h(u_\eta))_\eta \text{ is uniformly bounded in } L^1\left(0, T; (W^{1,q}(\Omega))'\right) \text{ for } q > d. \quad (5.57)$$

$$(h(u_\eta))_\eta \text{ and } (u_\eta)_\eta \text{ are relatively compact in } L^2(0, T; L^2(\Omega)). \quad (5.58)$$

Proof. The first part (5.53) is obtained in Lemma 5.10. Now, in order to obtain properties (5.54)–(5.56), we multiply the saturation equation (5.14) by $\beta(u_\eta)$, and get

$$E_1 + E_2 + E_3 + E_4 + E_5 + E_6 = 0. \quad (5.59)$$

where

$$\begin{aligned} E_1 &= \frac{d}{dt} \int_{\Omega} h(u_\eta) dx, \\ E_2 &= (r-1) \int_{\Omega} a(u_\eta) u_\eta^{r-2} \nabla u_\eta \cdot \nabla u_\eta dx = (r-1) \left\| \sqrt{u_\eta^{r-2} a(u_\eta)} \nabla u_\eta \right\|_{(L^2(\Omega))^d}^2, \\ E_3 &= \eta (r-1) \int_{\Omega} u_\eta^{r-2} |\nabla u_\eta|^2 dx = (r-1) \left\| \sqrt{\eta u_\eta^{r-2}} \nabla u_\eta \right\|_{(L^2(\Omega))^d}^2, \\ E_4 &= (r-1) \int_{\Omega} (g(u_\eta) - f(u_\eta)) u_\eta^{r-2} \nabla u_\eta \cdot \mathbf{V} dx + \int_{\Omega} g'(u_\eta) u_\eta^{r-1} \nabla u_\eta \cdot \mathbf{V} dx, \\ E_5 &= \int_{\Omega} a(u_\eta) u_\eta^{r-1} \nabla u_\eta \cdot \mathbf{V} dx, \quad E_6 = \eta \int_{\Omega} u_\eta^{r-1} \nabla u_\eta \cdot \mathbf{V} dx. \end{aligned}$$

Using assumptions (H2) and (H4b), the Cauchy-Schwarz and the weighted Young inequalities, we obtain the following estimate

$$\begin{aligned} |E_4| &\leq (r-1) C_{g'} \left(\int_{\Omega} |u_\eta^{r-1} \nabla u_\eta \cdot \mathbf{V}| dx + \int_{\Omega} |u_\eta^{r-1} \nabla u_\eta \cdot \mathbf{V}| dx \right) \\ &\leq C_{m_1, r} \|\mathbf{V}\|_{(L^2(\Omega))^d}^2 + \frac{m_1 (r-1)}{4} \|u_\eta^{r-1} \nabla u_\eta\|_{(L^2(\Omega))^d}^2. \end{aligned}$$

In the same manner, we get the estimates for the last two terms

$$\begin{aligned} |E_5| &\leq \int_{\Omega} |a(u_\eta) u_\eta^{r-2} \nabla u_\eta \cdot \mathbf{V}| dx \leq C_{a, r} \|\mathbf{V}\|_{(L^2(\Omega))^d}^2 + \frac{(r-1)}{4} \left\| \sqrt{u_\eta^{r-2} a(u_\eta)} \nabla u_\eta \right\|_{(L^2(\Omega))^d}^2, \\ |E_6| &\leq \int_{\Omega} |\eta u_\eta^{r-2} \nabla u_\eta \cdot \mathbf{V}| dx \leq \frac{1}{r-1} \|\mathbf{V}\|_{(L^2(\Omega))^d}^2 + \frac{(r-1)}{4} \left\| \eta^{\frac{1}{2}} u_\eta^{\frac{r-2}{2}} \nabla u_\eta \right\|_{(L^2(\Omega))^d}^2. \end{aligned}$$

Plugging the previous estimates into equation (5.59), this yields

$$\begin{aligned} \frac{d}{dt} \int_{\Omega} h(u_\eta) dx + (r-1) \int_{\Omega} a(u_\eta) u_\eta^{r-2} \nabla u_\eta \cdot \nabla u_\eta dx \\ + (r-1) \int_{\Omega} \eta u_\eta^{r-2} \nabla u_\eta \cdot \nabla u_\eta dx \leq C_{m_1, r, a}, \end{aligned} \quad (5.60)$$

where $C_{m_1, r, a}$ is a nonnegative constant independent of η .

Thanks to statement (5.53), then assumption (H4b) ensures that $m_1 u_\eta^r \leq m_1 u_\eta^{r_1} \leq a(u_\eta) \leq M_1 u_\eta^{r_1} \leq M_1 u_\eta^{r-2}$, and therefore we have the following inequalities

$$\begin{aligned} \|u_\eta^{r-1} \nabla u_\eta\|_{(L^2(\Omega))^d}^2 &\leq \frac{1}{m_1} \left\| \sqrt{u_\eta^{r-2} a(u_\eta)} \nabla u_\eta \right\|_{(L^2(\Omega))^d}^2, \text{ and} \\ \|a(u_\eta) \nabla u_\eta\|_{(L^2(\Omega))^d}^2 &\leq M_1 \left\| \sqrt{u_\eta^{r-2} a(u_\eta)} \nabla u_\eta \right\|_{(L^2(\Omega))^d}^2. \end{aligned}$$

Integrating inequality (5.60) with respect to the time $t \in (0, T)$, then estimates (5.54)–(5.56) hold. To show that $(\partial_t h(u_\eta))_\eta$ is uniformly bounded in $L^1\left(0, T; (W^{-1}(\Omega))'\right)$, we take a test function

$\chi \in L^2(0, T; H_0^1(\Omega)) \cap L^\infty(Q_T)$, multiply the saturation equation by $\beta(u_\eta)\chi$, and integrate the resulting equation over Q_T , one has

$$\begin{aligned}
& \int_0^T \langle \partial_t h(u_\eta), \chi \rangle dt = - \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla (\beta(u_\eta)\chi) dx dt \\
& - \eta \int_{Q_T} \nabla u_\eta \cdot \nabla (\beta(u_\eta)\chi) dx dt + \int_{Q_T} (f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla (\beta(u_\eta)\chi) dx dt \\
& - \int_{Q_T} g'(u_\eta) \nabla u_\eta \cdot \mathbf{V} \beta(u_\eta)\chi dx dt - \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \mathbf{V} \beta(u_\eta)\chi dx dt \\
& - \eta \int_{Q_T} \nabla u_\eta \cdot \mathbf{V} \beta(u_\eta)\chi dx dt.
\end{aligned} \tag{5.61}$$

Here, we give estimates on each integral on the right-hand side of equation (5.61) that we denote them I_i , $1 \leq i \leq 6$. To obtain the estimates, we use the Cauchy-Schwarz inequality. For the first term, we have

$$\begin{aligned}
|I_1| & \leq \int_{Q_T} |a(u_\eta) \nabla u_\eta \cdot ((r-1)u_\eta^{r-2}\nabla u_\eta\chi + u_\eta^{r-1}\nabla\chi)| dx dt \\
& \leq C_{r,a} \left(\int_{Q_T} |a(u_\eta) u_\eta^{r-2}\nabla u_\eta \cdot \nabla u_\eta\chi| dx dt + \int_{Q_T} |u_\eta^{r-1}\nabla u_\eta \cdot \nabla\chi| dx dt \right) \\
& \leq C_{r,a} \left\| \sqrt{a(u_\eta) u_\eta^{r-1}\nabla u_\eta} \right\|_{(L^2(Q_T))^d}^2 \|\chi\|_{L^\infty(Q_T)} \\
& \quad + \|u_\eta^{r-1}\nabla u_\eta\|_{(L^2(Q_T))^d} \|\nabla\chi\|_{(L^2(Q_T))^d}.
\end{aligned}$$

In the same manner, we have the estimate on the second term

$$\begin{aligned}
|I_2| & \leq \int_{Q_T} |\eta \nabla u_\eta \cdot ((r-1)u_\eta^{r-2}\nabla u_\eta\chi + u_\eta^{r-1}\nabla\chi)| dx dt \\
& \leq C_{r,a} \left\| \sqrt{\eta u_\eta^{r-2}\nabla u_\eta} \right\|_{(L^2(Q_T))^d} \|\chi\|_{L^\infty(Q_T)} \\
& \quad + \|u_\eta^{r-1}\nabla u_\eta\|_{(L^2(Q_T))^d} \|\nabla\chi\|_{(L^2(Q_T))^d}.
\end{aligned}$$

The third term is estimated, with the help of assumption (H2) and the Poincaré inequality, as follows

$$\begin{aligned}
|I_3| & \leq \int_{Q_T} |(f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot ((r-1)u_\eta^{r-2}\nabla u_\eta\chi + u_\eta^{r-1}\nabla\chi)| dx dt \\
& \leq C_{r,\Omega} \|\mathbf{V}\|_{(L^\infty(Q_T))^d}^2 \left(\|u_\eta^{r-1}\nabla u_\eta\|_{(L^2(Q_T))^d} + 1 \right) \|\nabla\chi\|_{(L^2(Q_T))^d}.
\end{aligned}$$

Similarly, we have

$$\begin{aligned}
|I_4| & \leq \int_{Q_T} |g'(u_\eta) u_\eta^{r-1}\nabla u_\eta \cdot \mathbf{V}\chi| dx dt \\
& \leq C_{g',\Omega} \left(\|u_\eta^{r-1}\nabla u_\eta\|_{(L^2(Q_T))^d} \|\mathbf{V}\|_{(L^\infty(Q_T))^d}^2 \right) \|\nabla\chi\|_{(L^2(Q_T))^d}.
\end{aligned}$$

Finally, the last two terms are estimated as

$$\begin{aligned}
|I_5 + I_6| & \leq \int_{Q_T} |a(u_\eta) u_\eta^{r-1}\nabla u_\eta \cdot \mathbf{V}\chi| dx dt + \int_{Q_T} |\eta u_\eta^{r-1}\nabla u_\eta \cdot \mathbf{V}\chi| dx dt \\
& \leq C_{a,\Omega} \left(\|u_\eta^{r-1}\nabla u_\eta\|_{(L^2(Q_T))^d} \|\mathbf{V}\|_{(L^\infty(Q_T))^d}^2 \right) \|\nabla\chi\|_{(L^2(Q_T))^d}.
\end{aligned}$$

Plugging the previous estimates into equation (5.61), one gets

$$|\langle \partial_t h(u_\eta), \chi \rangle| \leq \kappa \left(\|\chi\|_{L^\infty(Q_T)} + \|\chi\|_{L^2(0,T;H_0^1(\Omega))} \right).$$

One can conclude the proof of statement (5.57), using the boundedness of the sequences given in (5.53)–(5.55) and the embedding of the Sobolev space $W^{1,q}(\Omega) \subset H_0^1(\Omega) \cap L^\infty(\Omega)$ for $q > d$, and consequently, one has

$$L^\infty(0, T; W^{1,q}(\Omega)) \subset L^2(0, T; H^1(\Omega)) \cap L^\infty(0, T; L^2(\Omega)), \quad \forall q > d.$$

The last part (5.58) of the lemma is a direct consequence of the Aubin–Simon theorem. Indeed, one can prove that the sequence $(h(u_\eta))_\eta$ is uniformly bounded in $L^2(0, T; H_0^1(\Omega))$ since $\nabla h(u_\eta) = u_\eta^{r-1} \nabla u_\eta$. Then the sequence $(h(u_\eta))_\eta$ is lying into the Sobolev space

$$\mathcal{W} = \left\{ u_\eta; u_\eta \in L^2(0, T; H_0^1(\Omega)) \text{ and } \partial_t u_\eta \in L^1\left(0, T; (W^{1,q}(\Omega))'\right) \right\}.$$

Thanks to the Aubin–Lions lemma, \mathcal{W} is compactly embedded in $L^2(0, T; L^2(\Omega))$. Consequently, the sequence $(h(u_\eta))_\eta$ is relatively compact in $L^2(0, T; L^2(\Omega))$.

Since the differentiable function h is nondecreasing, then h^{-1} exists and it is continuous, then the sequence $(u_\eta)_\eta$ is relatively compact in $L^2(0, T; L^2(\Omega))$. The proof of lemma 5.13 is now complete. \square

Now, we are interested in an almost-everywhere convergence of the gradient of the saturation weighted by a degenerate function of the saturation. Specifically, we have the following lemma.

Lemma 5.14. *Let $q = 3r_1 + 2$, where r_1 is defined in assumption (H4b). Then, the sequence $(u_\eta^q a(u_\eta) \nabla u_\eta)_\eta$ is a Cauchy sequence in measure.*

Proof. Let us denote by A , B , and b the functions defined by

$$A(u) = \int_0^u a(\tau) \, d\tau, \quad B(u) = A^2(u), \quad b(u) = B'(u) = 2A(u) a(u). \quad (5.62)$$

For every $\mu > 0$ a nonnegative parameter, we consider the truncation function T_μ and its derivative Θ_μ defined by

$$T_\mu(u) = \min(\mu, \max(-\mu, u)), \quad \Theta_\mu(u) = \int_0^u T_\mu(\tau) \, d\tau, \quad \forall u \in \mathbb{R}. \quad (5.63)$$

We want to show that the sequence $(u_\eta^q a(u_\eta) \nabla u_\eta)_\eta$ is a Cauchy sequence in measure, this yields that, up to extract a subsequence, $u_\eta^q a(u_\eta) \nabla u_\eta \rightarrow u^q a(u) \nabla u$ almost everywhere in $\Omega(\mathbf{x}, t)$.

To do that, it suffices to prove that, for two sequences $(u_\eta)_\eta$ and $(u_{\eta'})_{\eta'}$ satisfying the saturation equation (5.14), we have

$$\text{meas} \left\{ \left| u_\eta^q \nabla A(u_\eta) - u_{\eta'}^q \nabla A(u_{\eta'}) \right| \geq \delta \right\} \leq \varepsilon, \quad \forall \varepsilon > 0. \quad (5.64)$$

We denote by $A_{\eta, \eta'} = A(u_\eta) - A(u_{\eta'})$ and $B_{\eta, \eta'} = B(u_\eta) - B(u_{\eta'})$.

We subtract the saturation equations satisfied by $(u_\eta)_\eta$ and $(u_{\eta'})_{\eta'}$, then we multiply by $\sigma_\eta = b(u_\eta) T_\mu(B_{\eta,\eta'})$ et $\sigma_{\eta'} = b(u_{\eta'}) T_\mu(B_{\eta,\eta'})$ respectively, one gets

$$\begin{aligned}
& \int_{\Omega} \theta_\mu(B_{\eta,\eta'}(t, x)) \, d\mathbf{x} + \int_{Q_t} (\nabla A(u_\eta) \cdot \nabla \sigma_\eta - \nabla A(u_{\eta'}) \cdot \nabla \sigma_{\eta'}) \, d\mathbf{x} \, dt \\
&= \int_{Q_t} ((f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla \sigma_\eta - (f(u_{\eta'}) - g(u_{\eta'})) \mathbf{V} \cdot \nabla \sigma_{\eta'}) \, d\mathbf{x} \, dt \\
&- \int_{Q_t} (g'(u_\eta) \nabla u_\eta \cdot \mathbf{V} \sigma_\eta - g'(u_{\eta'}) \nabla u_{\eta'} \cdot \mathbf{V} \sigma_{\eta'}) \, d\mathbf{x} \, dt \\
&- \int_{Q_t} \eta \nabla u_\eta \cdot \nabla \sigma_\eta \, d\mathbf{x} \, dt + \eta' \int_{Q_t} \nabla u_{\eta'} \cdot \nabla \sigma_{\eta'} \, d\mathbf{x} \, dt \\
&- \int_{Q_t} (\nabla A(u_\eta) \cdot \mathbf{V} \sigma_\eta - \nabla A(u_{\eta'}) \cdot \mathbf{V} \sigma_{\eta'}) \, d\mathbf{x} \, dt \\
&- \int_{Q_t} \eta \nabla u_\eta \cdot \mathbf{V} \sigma_\eta \, d\mathbf{x} \, dt + \eta' \int_{Q_t} \nabla u_{\eta'} \cdot \mathbf{V} \sigma_{\eta'} \, d\mathbf{x} \, dt.
\end{aligned} \tag{5.65}$$

Thanks to lemma 5.13 and to the assumption (H4b), one can deduce that the sequences

$$(\nabla A(u_\eta))_\eta, (\nabla B(u_\eta))_\eta, (\nabla b(u_\eta))_\eta \text{ are uniformly bounded in } (L^2(Q_T))^d, \tag{5.66}$$

Indeed, we have the following estimates

$$\begin{aligned}
\|\nabla A(u_\eta)\|_{(L^2(Q_T))^d}^2 &= \|a(u_\eta) \nabla u_\eta\|_{(L^2(Q_T))^d}^2 \leq C_{M_1, m_1, r, a}, \\
\|\nabla B(u_\eta)\|_{(L^2(Q_T))^d}^2 &= \|2A(u_\eta) \nabla A(u_\eta)\|_{(L^2(Q_T))^d}^2 \leq (2M_1)^2 \|\nabla A(u_\eta)\|_{(L^2(Q_T))^d}^2.
\end{aligned}$$

It remains to show that $(\nabla b(u_\eta))_\eta$ is uniformly bounded in $(L^2(Q_T))^d$. We have, from the definition of b , that $\nabla b(u_\eta) = 2a(u_\eta) \nabla A(u_\eta) + 2a'(u_\eta) A(u_\eta) \nabla u_\eta$. One can get the result using the following statement

$$|a'(u_\eta) A(u_\eta)| \leq M_1^2 r_1 u_\eta^{r_1-1} \frac{u_\eta^{r_1+1}}{r_1+1} \leq M_1^2 u_\eta^{2r_1} \leq \frac{M_1^2}{m_1} u_\eta^{r_1} a(u_\eta) \leq \frac{M_1^2}{m_1} a(u_\eta). \tag{5.67}$$

We denote by I_i , $i = 1, 7$, the integrals on the right-hand side of equation (5.65), and let $(\delta_\eta)_\eta$, $(\delta_{\eta'})_{\eta'}$, $(\mathbf{V}_\eta)_\eta$, and $(\mathbf{V}_{\eta'})_{\eta'}$ be the sequences given by

$$\delta_\eta = (f(u_\eta) - g(u_\eta)), \quad \delta_{\eta'} = (f(u_{\eta'}) - g(u_{\eta'})), \quad \mathbf{V}_\eta = \delta_\eta \mathbf{V}, \quad \mathbf{V}_{\eta'} = \delta_{\eta'} \mathbf{V}.$$

Using the Lebesgue theorem, we get

$$\|V_\eta - V_{\eta'}\|_{L^2(Q_T)} = \|(f(u_\eta) - g(u_\eta)) \mathbf{V} - (f(u_{\eta'}) - g(u_{\eta'})) \mathbf{V}\|_{L^2(Q_T)} \xrightarrow{\eta, \eta' \rightarrow 0} 0. \tag{5.68}$$

Now, we give estimates on each term on the right-hand side of equation (5.65). For the first term, we have

$$\begin{aligned}
& (f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla \sigma_\eta - (f(u_{\eta'}) - g(u_{\eta'})) \mathbf{V} \cdot \nabla \sigma_{\eta'} \\
&= (\mathbf{V}_\eta \cdot \nabla b(u_\eta) - \mathbf{V}_{\eta'} \cdot \nabla b(u_{\eta'})) T_\mu(B_{\eta,\eta'}) + (\mathbf{V}_\eta b(u_\eta) - \mathbf{V}_{\eta'} b(u_{\eta'})) \nabla T_\mu(B_{\eta,\eta'}) \\
&= (\mathbf{V}_\eta \cdot \nabla b(u_\eta) - \mathbf{V}_{\eta'} \cdot \nabla b(u_{\eta'})) T_\mu(B_{\eta,\eta'}) + (\mathbf{V}_\eta - \mathbf{V}_{\eta'}) b(u_\eta) \nabla T_\mu(B_{\eta,\eta'}) \\
&+ (b(u_\eta) - b(u_{\eta'})) \mathbf{V}_{\eta'} \cdot \nabla T_\mu(B_{\eta,\eta'}).
\end{aligned}$$

As a consequence,

$$\begin{aligned} |I_1| \leq & \left\| \left(\mathbf{V}_\eta \cdot \nabla b(u_\eta) - \mathbf{V}_{\eta'} \cdot \nabla b(u_{\eta'}) \right) T_\mu(B_{\eta, \eta'}) \right\|_{L^1(Q_T)} \\ & + \|b\|_{L^\infty(Q_T)}^2 \left\| \mathbf{V}_\eta - \mathbf{V}_{\eta'} \right\|_{(L^2(Q_T))^d} \left\| \nabla T_\mu(B_{\eta, \eta'}) \right\|_{(L^2(Q_T))^d} \\ & + \left\| (b(u_\eta) - b(u_{\eta'})) \mathbf{V}_{\eta'} \cdot \nabla T_\mu(B_{\eta, \eta'}) \right\|_{L^1(Q_T)}. \end{aligned} \quad (5.69)$$

The first term on the right-hand side of inequality (5.69) is estimated, using (5.66), as follows

$$\begin{aligned} & \left\| \left(\mathbf{V}_\eta \cdot \nabla b(u_\eta) - \mathbf{V}_{\eta'} \cdot \nabla b(u_{\eta'}) \right) T_\mu(B_{\eta, \eta'}) \right\|_{L^1(Q_T)} \\ & \leq \left\| \mathbf{V}_\eta \cdot \nabla b(u_\eta) T_\mu(B_{\eta, \eta'}) \right\|_{L^1(Q_T)} + \left\| \mathbf{V}_{\eta'} \cdot \nabla b(u_{\eta'}) T_\mu(B_{\eta, \eta'}) \right\|_{L^1(Q_T)} \\ & \leq C^2 \left\| \mathbf{V} \right\|_{(L^\infty(Q_T))^d}^2 \left\| T_\mu(B_{\eta, \eta'}) \right\|_{(L^2(Q_T))^d} \left(\left\| \nabla b(u_\eta) \right\|_{(L^2(Q_T))^d} + \left\| \nabla b(u_{\eta'}) \right\|_{(L^2(Q_T))^d} \right), \end{aligned}$$

where C is the nonnegative constant introduced in assumption (H2) independent of η .

Taking into account the uniform boundedness in $(L^2(Q_T))^d$ of the sequence $(\nabla b(u_\eta))_\eta$, and the following overestimate $|T_\mu(B_{\eta, \eta'})| \leq \mu$ as well as the convergence $B_{\eta, \eta'}$ to zero for almost every $(\mathbf{x}, t) \in Q_T$, one has

$$\left\| \left(\mathbf{V}_\eta \cdot \nabla b(u_\eta) - \mathbf{V}_{\eta'} \cdot \nabla b(u_{\eta'}) \right) T_\mu(B_{\eta, \eta'}) \right\|_{L^1(Q_T)} \xrightarrow{\eta, \eta' \rightarrow 0} 0.$$

It is easy to see, using the convergence (5.68), that the second term on the right-hand side of inequality (5.69) tends to zero as $\eta, \eta' \rightarrow 0$.

It remains to prove that the last term on the right-hand side of inequality (5.69) tends to zero as $\eta, \eta' \rightarrow 0$. Indeed, using the boundedness of the functions b (due to the boundedness of the dissipation function a and consequently to the function A), and thanks to the Lebesgue theorem, one has $\left\| (b(u_\eta) - b(u_{\eta'})) \mathbf{V}_{\eta'} \right\|_{(L^2(Q_T))^d}$ tends to zero as $\eta, \eta' \rightarrow 0$. One can conclude the result, using statement (5.66) and the boundedness of the sequence $(\nabla T_\mu(B_{\eta, \eta'}))_{\eta, \eta'}$.

For the second term of equation (5.65), we have using the definition of the function b that

$$\begin{aligned} & g'(u_\eta) \nabla u_\eta \cdot \mathbf{V} b(u_\eta) T_\mu(B_{\eta, \eta'}) - g'(u_{\eta'}) \nabla u_{\eta'} \cdot \mathbf{V} b(u_{\eta'}) T_\mu(B_{\eta, \eta'}) \\ & = (g'(u_\eta) \mathbf{V} \cdot \nabla B(u_\eta) - g'(u_{\eta'}) \mathbf{V} \cdot \nabla B(u_{\eta'})) T_\mu(B_{\eta, \eta'}). \end{aligned}$$

The Hölder inequality and Lebesgue's theorem ensure that

$$|I_2| \leq C_{g', \mathbf{V}} \left\| T_\mu(B_{\eta, \eta'}) \right\|_{L^2(Q_T)} \left\| \nabla B(u_\eta) \right\|_{(L^2(Q_T))^d} + \left\| \nabla B(u_\eta) \right\|_{(L^2(Q_T))^d} \xrightarrow{\eta, \eta' \rightarrow 0} 0.$$

For the third term I_3 of equation (5.65), we write

$$\begin{aligned} \nabla u_\eta \cdot \nabla (b(u_\eta) T_\mu(B_{\eta, \eta'})) &= 2 \nabla A(u_\eta) \cdot \nabla A(u_\eta) T_\mu(B_{\eta, \eta'}) \\ &+ 2 A(u_\eta) \cdot \nabla A(u_\eta) T_\mu(B_{\eta, \eta'}) + 2 A(u_\eta) a'(u_\eta) \nabla u_\eta \cdot \nabla u_\eta T_\mu(B_{\eta, \eta'}). \end{aligned}$$

Using statements (5.67)–(5.66), one can deduce with the help of the Hölder inequality, that $|I_4| \leq \kappa \eta$, for some constant $\kappa > 0$ independent of η and η' . Therefore, $|I_4| \xrightarrow{\eta, \eta' \rightarrow 0} 0$.

For the fifth term I_5 of equation (5.65), we write

$$\begin{aligned} I_5 &= \int_{Q_t} \nabla A(u_{\eta'}) \cdot \mathbf{V} b(u_{\eta'}) T_\mu(B_{\eta, \eta'}) - \nabla A(u_\eta) \cdot \mathbf{V} b(u_\eta) T_\mu(B_{\eta, \eta'}) \, d\mathbf{x} \, dt \\ &= \int_{Q_t} (a(u_{\eta'}) \mathbf{V} \cdot \nabla B(u_{\eta'}) - a(u_\eta) \mathbf{V} \cdot \nabla B(u_\eta)) T_\mu(B_{\eta, \eta'}) \, d\mathbf{x} \, dt. \end{aligned}$$

Using again Hölder's inequality and the Lebesgue theorem, one can conclude that

$$|I_5| \leq \left\| \left(a(u_\eta) \mathbf{V} \cdot \nabla B(u_\eta) - a(u_{\eta'}) \mathbf{V} \cdot \nabla B(u_{\eta'}) \right) T_\mu(B_{\eta, \eta'}) \right\|_{L^1(Q_T)} \xrightarrow{\eta, \eta' \rightarrow 0} 0.$$

Finally, for the last two terms of equation (5.65), we have

$$\int_{Q_t} \nabla u_\eta \cdot \mathbf{V} b(u_\eta) T_\mu(B_{\eta, \eta'}) \, d\mathbf{x} \, dt = 2 \int_{Q_t} \nabla A(u_\eta) \cdot \mathbf{V} A(u_\eta) T_\mu(B_{\eta, \eta'}) \, d\mathbf{x} \, dt.$$

As a consequence,

$$|I_6| \leq \eta \left(\|\nabla A(u_\eta)\|_{(L^2(Q_T))^d} \|\mathbf{V}\|_{(L^2(Q_T))^d} M_1 \mu \right) \leq \kappa \eta \xrightarrow{\eta, \eta' \rightarrow 0} 0.$$

Similarly, we prove that $|I_7| \leq \kappa \eta \xrightarrow{\eta, \eta' \rightarrow 0} 0$ for some constant $\kappa > 0$ independent of η and η' .

We denote by $W_\mu(\eta, \eta')$ the right-hand side of equation (5.65) and by $V(\mu)$ the firm term on the left-hand side of the same equation; from the estimations on the integrals I_i , $W_\mu(\eta, \eta')$ goes to zero as $\eta, \eta' \rightarrow 0$, for all $\mu > 0$. We also have $|V(\mu)| \leq |\Omega| \mu$, which goes to zero as $\mu \rightarrow 0$ and uniformly on η and η' . Therefore, we have the following result stemming from equation (5.65) and the aforementioned definitions

$$\begin{aligned} \int_{Q_T} (\nabla A(u_\eta) \cdot \nabla (b(u_\eta) T_\mu(B_{\eta, \eta'})) - \nabla A(u_{\eta'}) \cdot \nabla (b(u_{\eta'}) T_\mu(B_{\eta, \eta'}))) \, d\mathbf{x} \, dt \\ = W_\mu(\eta, \eta') + V(\mu). \end{aligned} \quad (5.70)$$

Here, we give the second step for the proof of statement (5.64). Let s the continuous function defined by

$$s(u) = \int_0^u b(z) A(z) a(z) \, dz, \quad \forall u \in \mathbb{R}. \quad (5.71)$$

Let us prove that

$$\text{meas}\{|\nabla s(u_\eta) - \nabla s(u_{\eta'})| \geq \delta\} \xrightarrow{\eta, \eta' \rightarrow 0} 0.$$

Remark that $\{|\nabla s(u_\eta) - \nabla s(u_{\eta'})| \geq \delta\} \subset \mathcal{A}_1 \cap \mathcal{A}_2 \cap \mathcal{A}_3 \cap \mathcal{A}_4$, where

$$\begin{aligned} \mathcal{A}_1 &= \{|\nabla A(u_\eta)| \geq k\}, \quad \mathcal{A}_2 = \{|\nabla A(u_{\eta'})| \geq k\}, \quad \mathcal{A}_3 = \{|B_{\eta, \eta'}| \geq k\}, \\ \mathcal{A}_4 &= \{|\nabla s(u_\eta) - \nabla s(u_{\eta'})| \geq \delta\} \cap \{|\nabla A(u_\eta)| \leq k\} \cap \{|\nabla A(u_{\eta'})| \leq k\} \cap \{|B_{\eta, \eta'}| \leq \mu\}. \end{aligned}$$

Thanks to the statement (5.66), and the continuous embedding of $L^2(Q_T)$ in $L^1(Q_T)$, we have

$$\begin{aligned} \text{meas}(\mathcal{A}_1) k &\leq \int_{\mathcal{A}_1} |\nabla A(u_\eta)| \, d\mathbf{x} \, dt \leq \int_{Q_T} |\nabla A(u_\eta)| \, d\mathbf{x} \, dt = \|\nabla A(u_\eta)\|_{(L^1(Q_T))^d} \\ &\leq C_\Omega \|\nabla A(u_\eta)\|_{(L^2(Q_T))^d} \leq C_{M_1, m_1, r, a, \Omega}. \end{aligned}$$

An analogous estimate holds for \mathcal{A}_2 . Therefore, by choosing k large enough, one gets $\text{meas}(\mathcal{A}_1) + \text{meas}(\mathcal{A}_2)$ is arbitrarily small. In the same manner, one gets

$$\text{mes}(\mathcal{A}_3) \leq \frac{1}{\mu} \|B_{\eta, \eta'}\|_{L^1(Q_T)},$$

which, for a fixed $\mu > 0$, tends to zero as $\eta, \eta' \rightarrow 0$.

It remains to show that $\text{meas}(\mathcal{A}_4)$ is small enough. Indeed, we have

$$\begin{aligned} |\nabla s(u_\eta) - \nabla s(u_{\eta'})|^2 &= |b(u_\eta) A(u_\eta) \nabla A(u_\eta) - b(u_{\eta'}) A(u_{\eta'}) \nabla A(u_{\eta'})|^2 \\ &= |b(u_\eta) A(u_\eta) \nabla A_{\eta, \eta'} + (b(u_\eta) A(u_\eta) - b(u_{\eta'}) A(u_{\eta'})) \nabla A(u_{\eta'})|^2, \end{aligned}$$

and therefore, one gets

$$\begin{aligned}
\delta \text{meas}(\mathcal{A}_4) &\leq \int_{\mathcal{A}_4} |\nabla s(u_\eta) - \nabla s(u_{\eta'})|^2 d\mathbf{x} dt \leq 2 \int_{\mathcal{A}_4} |b(u_\eta) A(u_\eta) \nabla A_{\eta, \eta'}|^2 d\mathbf{x} dt \\
&\quad + 2 \int_{\mathcal{A}_4} |b(u_\eta) A(u_\eta) - b(u_{\eta'}) A(u_{\eta'})|^2 |\nabla A(u_{\eta'})|^2 d\mathbf{x} dt \\
&\leq 4M_1^3 \int_{\mathcal{A}_4} b(u_\eta) (A(u_\eta) + A(u_{\eta'})) \nabla A_{\eta, \eta'} \cdot \nabla A_{\eta, \eta'} d\mathbf{x} dt \\
&\quad + 2k^2 \int_{Q_T} |b(u_\eta) A(u_\eta) - b(u_{\eta'}) A(u_{\eta'})|^2 d\mathbf{x} dt
\end{aligned}$$

The parameter k is chosen to be fixed and large enough; then the last term that we denote $W_k(\eta, \eta')$ goes to zero as $\eta, \eta' \rightarrow 0$. Consequently,

$$\begin{aligned}
\delta \text{meas}(\mathcal{A}_4) &\leq W_k(\eta, \eta') \\
&\quad + 4M_1^3 \int_{Q_T} b(u_\eta) (A(u_\eta) + A(u_{\eta'})) \nabla A_{\eta, \eta'} \cdot \nabla A_{\eta, \eta'} 1_{\{|B_{\eta, \eta'}| \leq \mu\}} d\mathbf{x} dt \quad (5.72)
\end{aligned}$$

We want to show that the last term on the left-hand side of inequality (5.72) is small enough. For that, we compute

$$\begin{aligned}
\nabla A_{\eta, \eta'} \cdot \nabla \sigma_\eta &= b(u_\eta) \nabla A_{\eta, \eta'} \cdot \nabla T_\mu(B_{\eta, \eta'}) + \nabla A_{\eta, \eta'} \cdot \nabla b(u_{\eta'}) T_\mu(B_{\eta, \eta'}) \\
&= b(u_\eta) \nabla A_{\eta, \eta'} \cdot \nabla B_{\eta, \eta'} 1_{\{|B_{\eta, \eta'}| \leq \mu\}} + \nabla A_{\eta, \eta'} \cdot \nabla b(u_{\eta'}) T_\mu(B_{\eta, \eta'}) \\
&= 2b(u_\eta) \nabla A_{\eta, \eta'} \cdot \nabla A(u_\eta) A(u_{\eta'}) 1_{\{|B_{\eta, \eta'}| \leq \mu\}} \\
&\quad - 2b(u_\eta) \nabla A_{\eta, \eta'} \cdot \nabla A(u_\eta) A(u_{\eta'}) 1_{\{|B_{\eta, \eta'}| \leq \mu\}} + \nabla A_{\eta, \eta'} \cdot \nabla b(u_{\eta'}) T_\mu(B_{\eta, \eta'}) \\
&= b(u_\eta) (A(u_\eta) + A(u_{\eta'})) \nabla A_{\eta, \eta'} \cdot \nabla A_{\eta, \eta'} 1_{\{|B_{\eta, \eta'}| \leq \mu\}} \\
&\quad + \left(b(u_\eta) A_{\eta, \eta'} \nabla (A(u_\eta) + A(u_{\eta'})) 1_{\{|B_{\eta, \eta'}| \leq \mu\}} + \nabla b(u_{\eta'}) T_\mu(B_{\eta, \eta'}) \right) \cdot \nabla A_{\eta, \eta'}. \quad (5.73)
\end{aligned}$$

We have the following estimates

$$\begin{aligned}
\|\nabla b(u_{\eta'}) T_\mu(B_{\eta, \eta'}) \cdot \nabla A_{\eta, \eta'}\|_{(L^1(Q_T))^d} \\
\leq \|T_\mu(B_{\eta, \eta'})\|_{L^\infty(Q_T)} \|\nabla A_{\eta, \eta'}\|_{(L^2(Q_T))^d} \|\nabla b(u_{\eta'})\|_{(L^2(Q_T))^d} \leq \kappa\mu,
\end{aligned}$$

and

$$\begin{aligned}
\|b(u_\eta) \nabla A_{\eta, \eta'} \cdot \nabla (A(u_\eta) + A(u_{\eta'})) A_{\eta, \eta'} 1_{\{|B_{\eta, \eta'}| \leq \mu\}}\|_{(L^1(Q_T))^d} \\
\leq 2\mu M_1^2 \|\nabla A_{\eta, \eta'}\|_{(L^2(Q_T))^d} \|\nabla (A(u_\eta) + A(u_{\eta'}))\|_{(L^2(Q_T))^d} \leq \kappa\mu.
\end{aligned}$$

Consequently, the two last terms of equation (5.73) tend to zero in $L^1(Q_T)$ as μ goes to zero and therefore, they can be included into the function $V(\mu)$. Now, we are interested in giving an

estimate on the left-hand side on equation (5.73). Indeed, we have

$$\begin{aligned}
\int_{Q_t} \nabla A_{\eta, \eta'} \cdot \nabla \sigma_\eta \, d\mathbf{x} \, dt &= \int_{Q_t} \nabla A_{\eta, \eta'} \cdot \nabla (b(u_\eta) T_\mu(B_{\eta, \eta'})) \, d\mathbf{x} \, dt \\
&= \int_{Q_T} (\nabla A(u_\eta) \cdot \nabla (b(u_\eta) T_\mu(B_{\eta, \eta'})) - \nabla A(u_{\eta'}) \cdot \nabla (b(u_{\eta'}) T_\mu(B_{\eta, \eta'}))) \, d\mathbf{x} \, dt \\
&\quad - \int_{Q_T} \nabla A(u_{\eta'}) \cdot \nabla (b(u_\eta) - b(u_{\eta'})) T_\mu(B_{\eta, \eta'}) \, d\mathbf{x} \, dt \\
&\quad - \int_{Q_T} (b(u_\eta) - b(u_{\eta'})) \nabla A(u_{\eta'}) \cdot \nabla B_{\eta, \eta'} 1_{\{|B_{\eta, \eta'}| \leq \mu\}} \, d\mathbf{x} \, dt.
\end{aligned}$$

It is easy to see, using the previous results, that

$$\int_{Q_t} \nabla A_{\eta, \eta'} \cdot \nabla \sigma_\eta \, d\mathbf{x} \, dt \leq W_\mu(\eta, \eta') + \kappa\mu \quad (5.74)$$

Finally, using estimates (5.73)–(5.74), then inequality (5.72) gives

$$\delta \text{mes}(\mathcal{A}_4) \leq W_\mu(\eta, \eta') + V(\mu) + W_k(\eta, \eta')$$

Using the above results, one can deduce that for all $\varepsilon > 0$, for all $\delta > 0$, there exists $\eta_0 > 0$ such that for all $\eta, \eta' \leq \eta_0$, we have

$$\text{meas}\{|\nabla s(u_\eta) - \nabla s(u_{\eta'})| \geq \delta\} \leq \varepsilon \quad (5.75)$$

Now, we can prove statement (5.64) with the help of inequality (5.75). Indeed, we have

$$u_\eta^q \nabla A(u_\eta) - u_{\eta'}^q \nabla A(u_{\eta'}) = u_\eta^q \nabla A_{\eta, \eta'} + (u_\eta^q - u_{\eta'}^q) \nabla A(u_{\eta'}).$$

Since $q = 3r_1 + 2$, then $u_\eta^q = u_\eta^{3r_1+2} \leq C_{r_1, m_1} b(u_\eta) A(u_\eta)$ where $C_{r_1, m_1} = \frac{(r_1 + 1)^2}{2m_1^3}$.

We write

$$\begin{aligned}
|u_\eta^q \nabla A_{\eta, \eta'}| &\leq C_{r_1, m_1} b(u_\eta) A(u_\eta) |\nabla A_{\eta, \eta'}| \\
&\leq C_{r_1, m_1} |\nabla s(u_\eta) - \nabla s(u_{\eta'})| + C_{r_1, m_1} |(b(u_\eta) A(u_\eta) - b(u_{\eta'}) A(u_{\eta'})) \nabla A(u_{\eta'})|.
\end{aligned}$$

Consequently,

$$\begin{aligned}
|u_\eta^q \nabla A(u_\eta) - u_{\eta'}^q \nabla A(u_{\eta'})| &\leq |(u_\eta^q - u_{\eta'}^q) \nabla A(u_{\eta'})| + C_{r_1, m_1} |\nabla s(u_\eta) - \nabla s(u_{\eta'})| \\
&\quad + C_{r_1, m_1} |(b(u_\eta) A(u_\eta) - b(u_{\eta'}) A(u_{\eta'})) \nabla A(u_{\eta'})|,
\end{aligned}$$

which converges to zero as $\eta, \eta' \rightarrow 0$. The result is due either to the convergence in $L^1(Q_T)$ for the first and the last terms on the right-hand side, or either by the help of (5.75). This ends the proof of lemma 5.14. \square

Thanks to lemma 5.13 and lemma 5.14, we have the following lemma

Lemma 5.15. *Let $q = 3r_1 + 2$, we have the following convergences as η goes to zero*

$$u_\eta \longrightarrow u \text{ strongly in } L^2(Q_T), \quad (5.76)$$

$$h_\theta(u_\eta) \longrightarrow h_\theta(u) \text{ strongly in } L^2(Q_T), \quad (5.77)$$

$$\sqrt{u_\eta^{r-2} a(u_\eta)} \nabla u_\eta \longrightarrow \sqrt{u^{r-2} a(u)} \nabla u \text{ weakly in } (L^2(Q_T))^d, \quad (5.78)$$

$$u_\eta^{r-1} \nabla u_\eta \longrightarrow u^{r-1} \nabla u \text{ weakly in } (L^2(Q_T))^d, \quad (5.79)$$

Furthermore,

$$u_\eta \longrightarrow u \text{ almost everywhere in } Q_T, \quad (5.80)$$

$$0 \leq u(\mathbf{x}, t) \leq 1 \text{ almost everywhere in } Q_T, \quad (5.81)$$

$$u_\eta^q a(u_\eta) \nabla u_\eta \longrightarrow u^q a(u) \nabla u \text{ almost everywhere in } (L^2(Q_T))^d, \quad (5.82)$$

Proof. To prove the strong convergence (5.76), we deduce from lemma (5.13), thanks to the Aubin–Lions lemma, that there exists a subsequence of $(u_\eta)_\eta$ such that

$$u_\eta \longrightarrow u \text{ strongly in } L^2(Q_T).$$

As a consequence, $h_\theta(u_\eta) \longrightarrow h_\theta(u)$ as $\eta \rightarrow 0$ strongly in $L^2(Q_T)$ since h_θ is a continuous function from $L^2(Q_T)$ to $L^2(Q_T)$.

To prove statement (5.78), we consider the continuous function $\gamma : L^2(Q_T) \longrightarrow L^2(Q_T)$ defined by $\gamma(u_\eta) = \int_0^{u_\eta} \sqrt{\tau^{r-2} a(\tau)} d\tau$. We have, $\gamma(u_\eta) \xrightarrow{\eta \rightarrow 0} \gamma(u)$ fortement dans $L^2(Q_T)$, and since the sequence $(\nabla \gamma(u_\eta))_\eta = (\sqrt{u_\eta a(u_\eta)} \nabla u_\eta)_\eta$ converges weakly in $(L^2(Q_T))^d$, then it is easy to see that

$$\nabla \gamma(u_\eta) = u_\eta^q a(u_\eta) \nabla u_\eta \xrightarrow{\eta \rightarrow 0} \nabla \gamma(u) = u^q a(u) \nabla u, \text{ weakly in } (L^2(Q_T))^d.$$

In the same manner, we prove the convergence result (5.79). The convergence (5.80) is a consequence of the strong convergence (5.76), while the statement (5.81) is obtained using the fact that $\{u \in L^2(Q_T), 0 \leq u(\mathbf{x}, t) \leq 1 \text{ for a.e. } (\mathbf{x}, t) \in Q_T\}$ is a closed subset of $L^2(Q_T)$.

Finally, the proof of the almost everywhere convergence (5.82) is a direct consequence of lemma 5.14. \square

To complete the proof of theorem 5.4, we pass to the limit as η goes to zero in the following weak formulation : $\forall \chi \in C^1([0, T]; H_0^1(\Omega))$ with $\chi(\cdot, T) = 0$,

$$\begin{aligned} & - \int_0^T h_\theta(u_\eta) \partial_t \chi \, d\mathbf{x} \, dt - \int_\Omega h_\theta(u_0) \chi(\mathbf{x}, 0) \, d\mathbf{x} + \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla \beta_\theta(u_\eta) \chi \, d\mathbf{x} \, dt \\ & + \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla \chi \beta_\theta(u_\eta) \, d\mathbf{x} \, dt + \eta \int_{Q_T} \nabla u_\eta \cdot \nabla \beta_\theta(u_\eta) \chi \, d\mathbf{x} \, dt \\ & + \eta \int_{Q_T} \nabla u_\eta \cdot \nabla \chi \beta_\theta(u_\eta) \, d\mathbf{x} \, dt - \int_{Q_T} (f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla \beta_\theta(u_\eta) \chi \, d\mathbf{x} \, dt \\ & + \int_{Q_T} g'(u_\eta) \mathbf{V} \cdot \nabla u_\eta \beta_\theta(u_\eta) \chi \, d\mathbf{x} \, dt - \int_{Q_T} (f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla \chi \beta_\theta(u_\eta) \, d\mathbf{x} \, dt \\ & + \int_{Q_T} a(u_\eta) \mathbf{V} \cdot \nabla u_\eta \beta_\theta(u_\eta) \chi \, d\mathbf{x} \, dt + \eta \int_{Q_T} \nabla u_\eta \cdot \mathbf{V} \beta_\theta(u_\eta) \chi \, d\mathbf{x} \, dt = 0. \end{aligned} \quad (5.83)$$

Indeed, we denote by $L_i, i = 1, \dots, 11$ the integral terms in equation (5.83).

From the definition of the continuous function β_θ , we have

$$\int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla \chi \beta_\theta(u_\eta) \, d\mathbf{x} \, dt = \int_{Q_T} u_\eta^{r-1} \nabla u_\eta \cdot u_\eta^\theta a(u_\eta) \nabla \chi \, d\mathbf{x} \, dt.$$

The sequence $(u_\eta^{r-1} \nabla u_\eta)_\eta$ converges weakly towards $u^{r-1} \nabla u$ in $(L^2(Q_T))^d$. Further, thanks to Lebesgue's theorem, the sequence $(u_\eta^\theta a(u_\eta) \nabla \chi)_\eta$ converges strongly towards $u^\theta a(u) \nabla \chi$ in

$(L^2(Q_T))^d$; this gives the convergences of terms L_4 , and L_{10} . In the same manner, we obtain the convergence of $L_8 + L_9$ towards

$$\int_{Q_T} g'(u) \mathbf{V} \cdot \nabla u \beta_\theta(u) \chi \, d\mathbf{x} \, dt - \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla \chi \beta_\theta(u) \, d\mathbf{x} \, dt.$$

Let us focus on the seventh term L_7 of equation (5.83). Since $\theta > 1$, then there exists $\theta_0 > 0$ such that $\theta = 1 + \theta_0$. Therefore, using Lebesgue's theorem and the weak convergence (5.79), one has

$$\begin{aligned} L_7 &= -(r-1+\theta) \int_{Q_T} u_\eta^{r-1} \nabla u_\eta \cdot \mathbf{V} u_\eta^{\theta_0} (f(u_\eta) - g(u_\eta)) \chi \, d\mathbf{x} \, dt \\ &\xrightarrow{\eta \rightarrow 0} \int_{Q_T} (f(u) - g(u)) \mathbf{V} \cdot \nabla \beta_\theta(u) \chi \, d\mathbf{x} \, dt. \end{aligned}$$

For the fifth term, we have

$$\begin{aligned} |L_5| &= \eta \left| \int_{Q_T} \nabla u_\eta \cdot \nabla \beta_\theta(u_\eta) \chi \, d\mathbf{x} \, dt \right| = \kappa \eta \int_{Q_T} \left| u_\eta^{r-2} u_\eta^\theta \chi \nabla u_\eta \cdot \nabla u_\eta \right| \, d\mathbf{x} \, dt \\ &\leq \kappa \eta \int_{Q_T} |u_\eta^{2r-2} \chi \nabla u_\eta \cdot \nabla u_\eta| \, d\mathbf{x} \, dt \leq \kappa \eta \|u_\eta^{r-1} \nabla u_\eta\|_{(L^2(Q_T))^d} \|\chi\|_{L^\infty(Q_T)}. \end{aligned}$$

As a consequence, $|L_5| \leq \kappa \eta \rightarrow 0$, as η goes to zero.

The convergence to zero for the sixth and the last terms, is similar to the that for L_7 . Indeed, we have

$$\begin{aligned} |L_6| &= \eta \left| \int_{Q_T} \nabla u_\eta \cdot \nabla \chi \beta_\theta(u_\eta) \, d\mathbf{x} \, dt \right| = \eta \left| \int_{Q_T} u_\eta^{r-1} \nabla u_\eta \cdot u_\eta^\theta \nabla \chi \, d\mathbf{x} \, dt \right| \\ &\leq \kappa \eta \|u_\eta^{r-1} \nabla u_\eta\|_{(L^2(Q_T))^d} \|\nabla \chi\|_{(L^2(Q_T))^d} \xrightarrow{\eta \rightarrow 0} 0. \end{aligned}$$

Now, we are interested in the third term L_3 of equation (5.83). One can remark that the sequence $(a(u_\eta) \nabla u_\eta \cdot \nabla \beta_\theta(u_\eta))_\eta \subset \mathbb{R}$ is nonnegative and

$$a(u_\eta) \nabla u_\eta \cdot \nabla \beta_\theta(u_\eta) = (r-1+\theta) u_\eta^{r-2+\theta} a(u_\eta) \nabla u_\eta \cdot \nabla u_\eta$$

converges almost everywhere, up to a subsequence, to $a(u) \nabla u \cdot \nabla \beta_\theta(u)$, since $r-2+\theta-2q-r_1 \geq 0$, i.e. $\theta \geq 7r_1 + 6 - r$.

Consider a nonnegative test function ($\chi \geq 0$); then the lemma of Fatou ensures that

$$\liminf_{\eta \rightarrow 0} \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla \beta_\theta(u_\eta) \chi \, d\mathbf{x} \, dt \geq \int_{Q_T} a(u) \nabla u \cdot \nabla \beta_\theta(u) \chi \, d\mathbf{x} \, dt,$$

then the limit solution u verifies inequality (5.6) given into definition 5.3. Finally, to obtain (5.7), we apply the Egorov theorem on the sequence $(a(u_\eta) \nabla u_\eta \cdot \nabla \beta_\theta(u_\eta))_\eta$ which converges almost everywhere. Indeed, we have

$$\begin{aligned} \forall \varepsilon > 0, \exists Q^\varepsilon \subset Q_T \text{ such that } \text{meas}(Q^\varepsilon) < \varepsilon, \text{ and} \\ a(u_\eta) \nabla u_\eta \cdot \nabla \beta_\theta(u_\eta) &\xrightarrow{\eta \rightarrow 0} a(u) \nabla u \cdot \nabla \beta_\theta(u) \text{ uniformly in } Q_T \setminus Q^\varepsilon. \end{aligned}$$

Now, we take a test function χ such that $\text{supp } \chi \subset ([0, T] \times \Omega) \setminus Q^\varepsilon$, then

$$\int_{Q_T \setminus Q^\varepsilon} a(u_\eta) \nabla u_\eta \cdot \nabla \beta_\theta(u_\eta) \chi \, d\mathbf{x} \, dt \xrightarrow{\eta \rightarrow 0} \int_{Q_T \setminus Q^\varepsilon} a(u) \nabla u \cdot \nabla \beta_\theta(u) \chi \, d\mathbf{x} \, dt.$$

The proof of theorem 5.4 is now accomplished. \square

5.6 Proof of theorem 5.6

The proof of theorem 5.6 is a combination of the two techniques used for theorem 5.2 and theorem 5.4.

Lemma 5.16. *Under the assumptions (H1) – (H3) and (H4c), assume that $G(u_0)$ belongs to $L^1(\Omega)$. Then the solutions to the saturation equation (5.14) verify*

- (i) $0 \leq u_\eta(\mathbf{x}, t) \leq 1$, for almost everywhere $(\mathbf{x}, t) \in Q_T$.
- (ii) The sequences $(\sqrt{\mu'(u_\eta)} a(u_\eta) \nabla u_\eta)_\eta$ and $(a(u_\eta) \nabla u_\eta)_\eta$ are uniformly bounded in $(L^2(Q_T))^d$.
- (iii) The sequences $(\sqrt{\eta \mu'(u_\eta)} \nabla u_\eta)_\eta$ and $(\nabla J(u_\eta))_\eta$ are uniformly bounded in $(L^2(Q_T))^d$.
- (iv) The sequence $(G(u_\eta))_\eta$ is uniformly bounded in $L^\infty(0, T; L^2(\Omega))$.
- (v) The sequence $(\partial_t J(u_\eta))_\eta$ is uniformly bounded in $L^1(0, T; W^{-1, q'}(\Omega))$.
- (vi) The sequences $(J(u_\eta))_\eta$ and $(u_\eta)_\eta$ are relatively compact in $L^2(0, T; L^2(\Omega))$.

Proof. The first part, (i), is obtained in section 5.3.2.

Now, we multiply the saturation equation (5.14) by $\mu(u_\eta)$ and integrate over Ω , this yields

$$\begin{aligned} & \frac{d}{dt} \int_{\Omega} G(u_\eta) d\mathbf{x} + \int_{\Omega} a(u_\eta) \mu'(u_\eta) |\nabla u_\eta|^2 d\mathbf{x} + \eta \int_{\Omega} \mu'(u_\eta) |\nabla u_\eta|^2 d\mathbf{x} \\ &= \int_{\Omega} (f(u_\eta) - g(u_\eta)) \mu'(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} - \int_{\Omega} g'(u_\eta) \mu(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} \quad (5.84) \\ & - \int_{\Omega} a(u_\eta) \mu(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} - \eta \int_{\Omega} \mu(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x}. \end{aligned}$$

We split the whole integral appearing in equation (5.84) into two parts, similarly as we do in section 5.4, we write $\int_{\Omega} = \int_{\Omega \cap \{u_\eta < u_*\}} + \int_{\Omega \cap \{u_\eta \geq u_*\}}$ and we denote $\Omega_1 = \Omega \cap \{u_\eta < u_*\}$ and $\Omega_2 = \Omega \cap \{u_\eta \geq u_*\}$.

• Into the region Ω_1 ; and following the same analysis as the one developed into section 5.4, we obtain the following estimates

$$\begin{aligned} & \left| \int_{\Omega_1} (f(u_\eta) - g(u_\eta)) \mu'(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} - \int_{\Omega_1} g'(u_\eta) \mu(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} \right| \\ & \leq \frac{1}{4} \left\| \sqrt{\mu'(u_\eta)} a(u_\eta) \nabla u_\eta \right\|_{(L^2(\Omega_1))^d}^2 + 2 \frac{1 + (C(r-1))^2}{m_1(r-1)} \|\mathbf{V}\|_{(L^2(\Omega_1))^d}^2, \\ & \left| \int_{\Omega_1} a(u_\eta) \mu(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} \right| \leq \frac{1}{4} \left\| \sqrt{\mu'(u_\eta)} a(u_\eta) \nabla u_\eta \right\|_{(L^2(\Omega_1))^d}^2 + \frac{\|a\|_\infty^2}{r-1} \|\mathbf{V}\|_{(L^2(\Omega_1))^d}^2, \\ & \left| \eta \int_{\Omega_1} \mu(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} \right| \leq \frac{1}{2} \left\| \sqrt{\eta \mu'(u_\eta)} \nabla u_\eta \right\|_{(L^2(\Omega_1))^d}^2 + \frac{1}{2(r-1)} \|\mathbf{V}\|_{(L^2(\Omega_1))^d}^2. \end{aligned}$$

• Into region Ω_2 , we have

$$\int_{\Omega_2} ((f(u_\eta) - g(u_\eta)) \mu'(u_\eta) - g'(u_\eta) \mu(u_\eta)) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} = 0.$$

Furthermore, we have the following estimates

$$\begin{aligned} & \left| \int_{\Omega_2} a(u_\eta) \mu(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} \right| \\ & \leq \frac{1}{2} \left\| \sqrt{\mu'(u_\eta) a(u_\eta)} \nabla u_\eta \right\|_{(L^2(\Omega_2))^d}^2 + \frac{k_1 \|\sqrt{\mu} \mu' a\|_\infty^2}{4} \|\mathbf{V}\|_{(L^2(\Omega_2))^d}^2, \\ & \left| \eta \int_{\Omega_2} \mu(u_\eta) \nabla u_\eta \cdot \mathbf{V} d\mathbf{x} \right| \leq \frac{1}{2} \left\| \sqrt{\eta \mu'(u_\eta)} \nabla u_\eta \right\|_{(L^2(\Omega_2))^d}^2 + \frac{k_1 \|\sqrt{\mu} \mu'\|_\infty^2}{4} \|\mathbf{V}\|_{(L^2(\Omega_2))^d}^2. \end{aligned}$$

Plugging the previous estimates into equation (5.84), one has

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_{\Omega} G(u_\eta) d\mathbf{x} + \left\| \sqrt{\mu'(u_\eta) a(u_\eta)} \nabla u_\eta \right\|_{(L^2(\Omega))^d}^2 + \left\| \sqrt{\eta \mu'(u_\eta)} \nabla u_\eta \right\|_{(L^2(\Omega))^d}^2 \\ & \leq \kappa \|\mathbf{V}\|_{(L^2(\Omega))^d}^2 \end{aligned} \quad (5.85)$$

Now, we integrate inequality (5.85) with respect to the time over $(0, t)$, $t \in (0, T)$, one deduces that the sequences $(\sqrt{\mu'(u_\eta) a(u_\eta)} \nabla u_\eta)_\eta$ and $(\sqrt{\eta \mu'(u_\eta)} \nabla u_\eta)_\eta$ are uniformly bounded in $(L^2(Q_T))^d$, and that $(G(u_\eta))_\eta$ is uniformly bounded in $L^\infty(0, T; L^2(\Omega))$.

Let us prove that $(a(u_\eta) \nabla u_\eta)_\eta$ and $(\nabla J(u_\eta))_\eta$ are uniformly bounded in $(L^2(Q_T))^d$. Indeed, for all $0 \leq u \leq u_*$, we have

$$a(u) \mu'(u) \geq m_1 (r-1) u^{r_1} u^{r-2} \geq m_1 (r-1) u^r u^{r-2} \geq m_1 (r-1) u^{2r-2} \geq c_1 j^2(u),$$

and for all $u \geq u_*$, we have

$$\begin{aligned} a(u) \mu'(u) & \geq m_1 (1-u)^{r_2} \mu(u) g'(u) (f(u) - g(u))^{-1} \geq c_1 (1-u)^{r_2-2} \\ & \geq \frac{c_1}{(1-u_*)^{2-r'}} \left((1-u_*)^{1-\frac{r'}{2}} \right)^2 \left((1-u)^{\frac{r'}{2}-1} \right)^2 \geq \kappa j^2(u), \end{aligned}$$

Therefore, $\|\nabla J(u_\eta)\|_{(L^2(Q_T))^d}^2 \leq \kappa \left\| \mu'^{\frac{1}{2}}(u_\eta) a^{\frac{1}{2}}(u_\eta) \nabla u_\eta \right\|_{(L^2(Q_T))^d}^2 \leq \kappa$.

For the sequence $(a(u_\eta) \nabla u_\eta)_\eta$. It is easy to see that

$$\begin{aligned} a(u_\eta) & \leq M_1 u_\eta^{r_1} \leq M_1 u_\eta^{r-2} \leq \frac{M_1}{r-1} \mu'(u_\eta), & \text{if } 0 \leq u_\eta \leq u_*, \\ a(u_\eta) & \leq M_1 (1-u_\eta)^{r_2} \leq (1-u_\eta)^{-2} \leq \kappa \mu'(u_\eta), & \text{if } u_* \leq u_\eta \leq 1. \end{aligned}$$

As a consequence, the sequence $(a(u_\eta) \nabla u_\eta)_\eta$ is uniformly bounded in $(L^2(Q_T))^d$.

Let us now focus on the fourth part (iv), we want to prove that

$$(\partial_t J(u_\eta))_\eta \text{ is uniformly bounded in } L^2\left(0, T; (H^1(\Omega))'\right) + L^1(Q_T).$$

We take a test function $\chi \in L^2(0, T; H_0^1(\Omega)) \cap L^\infty(Q_T)$ and multiply the saturation equation

(5.14) by $j(u_\eta)\chi$, this yields

$$\begin{aligned}
\langle \partial_t J(u_\eta), \chi \rangle &= - \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla (j(u_\eta)\chi) \, d\mathbf{x} \, dt \\
&\quad - \eta \int_{Q_T} \nabla u_\eta \cdot \nabla (j(u_\eta)\chi) \, d\mathbf{x} \, dt + \int_{Q_T} (f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla (j(u_\eta)\chi) \, d\mathbf{x} \, dt \\
&\quad - \int_{Q_T} g'(u_\eta) \nabla u_\eta \cdot \mathbf{V} j(u_\eta)\chi \, d\mathbf{x} \, dt - \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \mathbf{V} j(u_\eta)\chi \, d\mathbf{x} \, dt \\
&\quad - \eta \int_{Q_T} \nabla u_\eta \cdot \mathbf{V} j(u_\eta)\chi \, d\mathbf{x} \, dt.
\end{aligned} \tag{5.86}$$

We will give estimates on each integral on the right-hand side of equation (5.86).

Into region $Q_T \cap \{u_\eta < u_*\}$, we have $j(u_\eta) = \mu(u_\eta)$, thus we obtain the same estimates on the integrals obtained within the proof of theorem 5.4. It remains to estimate the terms of the form

$\int_{Q_T \cap \{u_\eta \geq u_*\}}$ that we denote them by $\{L_i\}_{1 \leq i \leq 6}$ respectively.

For the first term L_1 , we have

$$\begin{aligned}
|L_1| &\leq \int_{Q_T \cap \{u_\eta \geq u_*\}} |a(u_\eta) j'(u_\eta) \nabla u_\eta \cdot \nabla u_\eta \chi| + |a(u_\eta) j(u_\eta) \nabla u_\eta \cdot \nabla \chi| \, d\mathbf{x} \, dt \\
&\leq \left\| \sqrt{j'(u_\eta) a(u_\eta) \nabla u_\eta} \right\|_{(L^2(Q_T))^d}^2 \|\chi\|_{L^\infty(Q_T)} + \|a(u_\eta) \nabla J(u_\eta)\|_{(L^2(Q_T))^d} \|\nabla \chi\|_{(L^2(Q_T))^d}.
\end{aligned}$$

On the other hand, using the definition of μ and j , we have, for all $u_* \leq u \leq 1$, that

$$|j'(u_\eta)| = \left(\frac{r'}{2} - 1 \right) \beta(u_*) (1 - u_*)^{1 - \frac{r'}{2}} (1 - u_\eta)^{\frac{r'}{2} - 2} \leq C_{u_*, r'} (1 - u_\eta)^{-2} \leq \mu'(u_\eta),$$

thus, thanks to parts (i) and (ii), one deduces that $|L_1| \leq \kappa \left(\|\chi\|_{L^\infty(Q_T)} + \|\nabla \chi\|_{(L^2(Q_T))^d} \right)$.

In the same manner, we obtain the estimates on the remaining terms except the estimate on L_3 . Indeed, using the assumption (H2) on f and g , one has

$$f(u_\eta) - g(u_\eta) \leq C_{g'} (1 - u_\eta), \quad \forall u_* \leq u_\eta \leq 1,$$

and therefore, we obtain the following estimates

$$\begin{aligned}
\left| \int_{Q_T} (f(u_\eta) - g(u_\eta)) j'(u_\eta) \nabla u_\eta \cdot \mathbf{V} \chi \, d\mathbf{x} \, dt \right| &\leq C_{g'} \int_{Q_T} |j(u_\eta) \nabla u_\eta \cdot \mathbf{V} \chi| \, d\mathbf{x} \, dt \\
&\leq C_{g'} \|\nabla J(u_\eta)\|_{(L^2(Q_T))^d} \|\mathbf{V}\|_{(L^2(Q_T))^d} \|\chi\|_{L^\infty(Q_T)}. \\
\left| \int_{Q_T} (f(u_\eta) - g(u_\eta)) j(u_\eta) \mathbf{V} \cdot \nabla \chi \, d\mathbf{x} \, dt \right| &\leq C_j \int_{Q_T} |\mathbf{V} \cdot \nabla \chi| \, d\mathbf{x} \, dt \\
&\leq C \|\mathbf{V}\|_{(L^2(Q_T))^d} \|\chi\|_{(L^2(Q_T))^d}.
\end{aligned}$$

Plugging the previous estimates into equation (5.86), one gets

$$|\langle \partial_t J(u_\eta), \chi \rangle| \leq \kappa \left(\|\chi\|_{L^\infty(Q_T)} + \|\chi\|_{L^2(0, T; H_0^1(\Omega))} \right).$$

One can conclude the proof of part (iii), using the embedding of the Sobolev space $W^{1, q}(\Omega) \subset H_0^1(\Omega) \cap L^\infty(\Omega)$ for $q > d$, and consequently, one has

$$L^\infty(0, T; W^{1, q}(\Omega)) \subset L^2(0, T; H^1(\Omega)) \cap L^\infty(0, T; L^2(\Omega)), \quad \forall q > d.$$

To complete the proof of the lemma, we remark that the sequence $(J(u_\eta))_\eta$ is lying into the Sobolev space

$$\mathcal{W} = \left\{ J(u_\eta); J(u_\eta) \in L^2(0, T; H_0^1(\Omega)) \text{ and } \partial_t J(u_\eta) \in L^\infty(0, T; W^{-1, q'}(\Omega)) \right\}.$$

Thanks to the Aubin–Simon theorem, \mathcal{W} is compactly embedded in $L^2(Q_T)$, and the sequence $(J(u_\eta))_\eta$ is relatively compact in $L^2(0, T; L^2(\Omega))$.

Since the differentiable function J is nondecreasing, then J^{-1} exists and it is continuous, then the sequence $(u_\eta)_\eta$ is relatively compact in $L^2(0, T; L^2(\Omega))$. The proof of lemma 5.16 is now accomplished. \square

Lemma 5.17. *Let $q_1 = 3r_1 + 2$ and $q_2 = 3r_2 + 2$, where r_1 and r_2 are given in assumption (H4c). The sequences $(1_{\{u_\eta \leq u_*\}} u_\eta^{q_1} a(u_\eta) \nabla u_\eta)_\eta$ and $(1_{\{u_\eta \geq u_*\}} (1 - u_\eta)^{q_2} a(u_\eta) \nabla u_\eta)_\eta$ are two Cauchy sequences in measure.*

Proof. The proof follows the same guidelines as the proof of lemma 5.14. We omit it for the sake of brevity. \square

To conclude the proof of theorem 5.6, we deduce from lemma 5.16 and lemma 5.17, that we can extract a subsequence such that we have the following convergences

$$\begin{aligned} 0 &\leq u(\mathbf{x}, t) \leq 1 \text{ a.e. in } Q_T, \\ u_\eta &\longrightarrow u \text{ strongly in } L^2(Q_T) \text{ and a.e. in } Q_T, \\ J(u_\eta) &\longrightarrow J(u) \text{ strongly } L^2(Q_T), \\ J(u_\eta) &\longrightarrow J(u) \text{ weakly in } L^2(0, T; H_0^1(\Omega)), \\ \sqrt{a(u_\eta) \mu'(u_\eta)} \nabla u_\eta &\longrightarrow \sqrt{a(u) \mu'(u)} \nabla u \text{ weakly in } (L^2(Q_T))^d, \\ 1_{\{u_\eta \leq u_*\}} u_\eta^{q_1} a(u_\eta) \nabla u_\eta &\longrightarrow 1_{\{u_\eta \leq u_*\}} u^{q_1} a(u) \nabla u \text{ a.e. in } Q_T, \\ 1_{\{u_\eta \geq u_*\}} (1 - u_\eta)^{q_2} a(u_\eta) \nabla u_\eta &\longrightarrow 1_{\{u_\eta \geq u_*\}} (1 - u)^{q_2} a(u) \nabla u \text{ a.e. in } Q_T. \end{aligned} \tag{5.87}$$

We consider the following weak formulation

$$\begin{aligned} & - \int_{Q_T} J_{\theta, \lambda}(u_\eta) \partial_t \chi \, d\mathbf{x} \, dt - \int_{\Omega} J_{\theta, \lambda}(u_0(\mathbf{x})) \chi(\mathbf{x}, 0) \, d\mathbf{x} \\ & + \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla j_{\theta, \lambda}(u_\eta) \chi \, d\mathbf{x} \, dt + \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla \chi j_{\theta, \lambda}(u_\eta) \, d\mathbf{x} \, dt \\ & + \eta \int_{Q_T} \nabla u_\eta \cdot \nabla j_{\theta, \lambda}(u_\eta) \chi \, d\mathbf{x} \, dt + \eta \int_{Q_T} \nabla u_\eta \cdot \nabla \chi j_{\theta, \lambda}(u_\eta) \, d\mathbf{x} \, dt \\ & - \int_{Q_T} (f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla j_{\theta, \lambda}(u_\eta) \chi \, d\mathbf{x} \, dt + \int_{Q_T} g'(u_\eta) \mathbf{V} \cdot \nabla u_\eta j_{\theta, \lambda}(u_\eta) \chi \, d\mathbf{x} \, dt \\ & - \int_{Q_T} (f(u_\eta) - g(u_\eta)) \mathbf{V} \cdot \nabla \chi j_{\theta, \lambda}(u_\eta) \, d\mathbf{x} \, dt + \int_{Q_T} a(u_\eta) \mathbf{V} \cdot \nabla u_\eta j_{\theta, \lambda}(u_\eta) \chi \, d\mathbf{x} \, dt \\ & + \eta \int_{Q_T} \nabla u_\eta \cdot \mathbf{V} j_{\theta, \lambda}(u_\eta) \chi \, d\mathbf{x} \, dt = 0, \quad \forall \chi \in \mathcal{C}^1([0, T]; H_0^1(\Omega)) \text{ with } \chi(T, \cdot) = 0 \end{aligned} \tag{5.88}$$

By splitting these integrals into two sub integrals, then we have the same convergence results obtained within the proof of theorem 5.4 for the integrals of the form $\int_{Q_T \cap \{u_\eta \leq u_*\}}$ in (5.88).

Now, we show the convergence for the remaining terms, we have

$$\int_{\{u_\eta \geq u_*\}} a(u_\eta) \nabla u_\eta \cdot \nabla \chi j_{\theta, \lambda}(u_\eta) \, d\mathbf{x} \, dt \xrightarrow{\eta \rightarrow 0} \int_{\{u \geq u_*\}} a(u) \nabla u \cdot \nabla \chi j_{\theta, \lambda}(u) \, d\mathbf{x} \, dt,$$

Indeed, on the one hand, the sequence $(a(u_\eta) \nabla u_\eta)_\eta$ converges weakly in $(L^2(Q_T))^d$ towards $a(u) \nabla u$. On the other hand, the sequence $(\nabla \chi_{j_{\theta,\lambda}}(u_\eta))_\eta$ converges strongly in $(L^2(Q_T))^d$ towards $\nabla \chi_{j_{\theta,\lambda}}(u)$. Furthermore, one obtains that

$$\eta \int_{\{u_\eta \geq u_*\}} \nabla u_\eta \cdot \nabla j_{\theta,\lambda}(u_\eta) \chi \, d\mathbf{x} \, dt \xrightarrow{\eta \rightarrow 0} 0. \quad (5.89)$$

Indeed,

$$\begin{aligned} \eta \left| \int_{\{u_\eta \geq u_*\}} \nabla u_\eta \cdot \nabla j_{\theta,\lambda}(u_\eta) \chi \, d\mathbf{x} \, dt \right| &= \kappa \eta \int_{Q_T} \left| (1 - u_\eta)^{\frac{r'}{2} - 2 + \lambda} \chi \nabla u_\eta \cdot \nabla u_\eta \right| \, d\mathbf{x} \, dt \\ &\leq \kappa \eta \int_{Q_T} \left| (1 - u_\eta)^{2r_2} \chi \nabla u_\eta \cdot \nabla u_\eta \right| \, d\mathbf{x} \, dt \leq \kappa \eta \|a(u_\eta) \nabla u_\eta\|_{(L^2(Q_T))^d} \|\chi\|_{L^\infty(Q_T)}. \end{aligned}$$

Then the convergence result (5.89) holds.

Now, we are interested in the third term of equation (5.88). We have $(a(u_\eta) \nabla u_\eta \cdot \nabla j_{\theta,\lambda}(u_\eta))_\eta$ is a nonnegative sequence and

$$1_{\{u_\eta \geq u_*\}} a(u_\eta) \nabla u_\eta \cdot \nabla j_{\theta,\lambda}(u_\eta) = 1_{\{u_\eta \geq u_*\}} c(u_*) (1 - u_\eta)^{\frac{r'}{2} - 2 + \lambda} a(u_\eta) \nabla u_\eta \cdot \nabla u_\eta$$

which converges almost everywhere, up to a subsequence, to $1_{\{u \geq u_*\}} a(u) \nabla u \cdot \nabla j_{\theta,\lambda}(u)$, since $\frac{r'}{2} - 2 + \lambda - 2q_2 - r_2 \geq 0$, i.e. $\lambda \geq 7r_2 + 6 - \frac{r'}{2}$.

Consider a nonnegative test function $(\chi \geq 0)$; then the Fatou's lemma ensures that

$$\liminf_{\eta \rightarrow 0} \int_{Q_T} a(u_\eta) \nabla u_\eta \cdot \nabla j_{\theta,\lambda}(u_\eta) \chi \, d\mathbf{x} \, dt \geq \int_{Q_T} a(u) \nabla u \cdot \nabla j_{\theta,\lambda}(u) \chi \, d\mathbf{x} \, dt,$$

then the limit solution u verifies inequality (5.12) into definition 5.5. Finally, to obtain (5.13), we apply the Egorov theorem on the sequence $(a(u_\eta) \nabla u_\eta \cdot \nabla j_{\theta,\lambda}(u_\eta))_\eta$ which converges almost everywhere. Indeed, we have

$$\begin{aligned} \forall \varepsilon > 0, \exists Q^\varepsilon \subset Q_T \text{ tel que } \text{mes}(Q^\varepsilon) < \varepsilon, \text{ and} \\ a(u_\eta) \nabla u_\eta \cdot \nabla j_{\theta,\lambda}(u_\eta) &\xrightarrow{\eta \rightarrow 0} a(u) \nabla u \cdot \nabla j_{\theta,\lambda}(u) \text{ uniformly in } Q_T \setminus Q^\varepsilon. \end{aligned}$$

Now, we take a nonnegative test function χ such that $\text{supp} \chi \subset ([0, T] \times \Omega) \setminus Q^\varepsilon$, then

$$\int_{Q_T \setminus Q^\varepsilon} a(u_\eta) \nabla u_\eta \cdot \nabla j_{\theta,\lambda}(u_\eta) \chi \, d\mathbf{x} \, dt \xrightarrow{\eta \rightarrow 0} \int_{Q_T \setminus Q^\varepsilon} a(u) \nabla u \cdot \nabla j_{\theta,\lambda}(u) \chi \, d\mathbf{x} \, dt.$$

This ends the proof of theorem 5.6. □



Technical Lemmas

We give in this appendix some technical results stemming from finite element discretizations. Let Ω be an open bounded polygonal subset of \mathbb{R}^2 , and \mathcal{T} be a conforming triangular mesh of Ω . For all triangle $T \in \mathcal{T}$, we denote by $h_T = \text{diam}(T)$, and by $h = \sup_{T \in \mathcal{T}} h_T$. We consider the set \mathcal{V} of the vertices of the triangles, and $(\mathbf{x}_K)_{K \in \mathcal{V}}$ their coordinates in Ω , and the usual \mathbb{P}_1 finite element space $\mathcal{H}_{\mathcal{T}}$ defined by

$$\mathcal{H}_{\mathcal{T}} = \{\phi \in C^0(\overline{\Omega}) ; \phi|_T \in \mathbb{P}_1(\mathbb{R}), \quad \forall T \in \mathcal{T}\}.$$

Let $u_{\mathcal{T}} \in \mathcal{H}_{\mathcal{T}}$, then consider the piecewise constant functions $\bar{u}_{\mathcal{T}}, \underline{u}_{\mathcal{T}} : \Omega \rightarrow \mathbb{R}$ defined by

$$\begin{aligned} \bar{u}_{\mathcal{T}}(\mathbf{x}) &= \bar{u}_T = \sup_{\mathbf{x} \in T} u_{\mathcal{T}}(\mathbf{x}), \quad \text{if } \mathbf{x} \in T \in \mathcal{T}, \\ \underline{u}_{\mathcal{T}}(\mathbf{x}) &= \underline{u}_T = \inf_{\mathbf{x} \in T} u_{\mathcal{T}}(\mathbf{x}), \quad \text{if } \mathbf{x} \in T \in \mathcal{T}. \end{aligned}$$

Lemma A.1. *There exists an absolute constant $c > 0$ such that*

$$\int_{\Omega} |\bar{u}_{\mathcal{T}}(\mathbf{x}) - \underline{u}_{\mathcal{T}}(\mathbf{x})|^2 d\mathbf{x} \leq ch^2 \int_{\Omega} |\nabla u_{\mathcal{T}}(\mathbf{x})|^2 d\mathbf{x},$$

where $c = \frac{243}{2\pi^2}$.

Proof. Let $T \in \mathcal{T}$ be a triangle whose vertices are located at $\mathbf{x}_K, \mathbf{x}_L$, and \mathbf{x}_M . We assume without loss of generality that

$$\underline{u}_T = u_M \leq u_L \leq u_K = \bar{u}_T \tag{A.1}$$

We denote by $\omega_{K,L}^T$ the triangle whose vertices are $\mathbf{x}_K, \mathbf{x}_{KL}$, and \mathbf{x}_T , where \mathbf{x}_T and \mathbf{x}_{KL} represent respectively, the center of gravity and the midpoint of $[\mathbf{x}_K, \mathbf{x}_L]$ of the triangle T . Specifically, we have

$$\mathbf{x}_T = \frac{1}{3}(\mathbf{x}_K + \mathbf{x}_L + \mathbf{x}_M), \quad \mathbf{x}_{KL} = \frac{1}{2}(\mathbf{x}_K + \mathbf{x}_L).$$

Note that

$$|\omega_{K,L}^T| = \frac{|T|}{6}, \tag{A.2}$$

where $|\omega_{K,L}^T|$ (respectively $|T|$) denotes the Lebesgue measure of $\omega_{K,L}^T$ (respectively T). Since $u_T \in \mathcal{H}_T$ is a linear function of \mathbf{x} on T , then one has

$$u_T := \frac{1}{|T|} \int_T u_T(\mathbf{x}) \, d\mathbf{x} = u(\mathbf{x}_T) = \frac{1}{3} (u_K + u_L + u_M),$$

and in view of property (A.1), one has

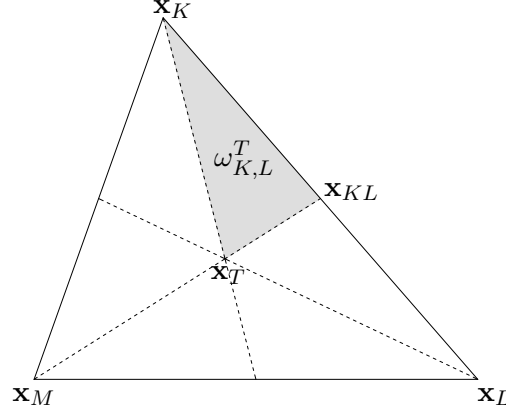


FIGURE A.1 – The triangle $T \in \mathcal{T}$ and the sub-triangle $w_{K,L}^T$.

$$u_T(\mathbf{x}_{KL}) = \frac{1}{2} (u_K + u_L) \geq u_T.$$

As a consequence,

$$u_T(\mathbf{x}) \geq u_T, \quad \forall \mathbf{x} \in \omega_{K,L}^T. \quad (\text{A.3})$$

Therefore, one has

$$\begin{aligned} \int_T |u_T(\mathbf{x}) - u_T|^2 \, d\mathbf{x} &\geq \int_{\omega_{K,L}^T} (u_T(\mathbf{x}) - u_T)^2 \, d\mathbf{x} = \frac{|\omega_{K,L}^T|}{9} (u_K + u_{KL} - 2u_T)^2 \\ &= \frac{|\omega_{K,L}^T|}{9} \left(u_K + \frac{u_K + u_L}{2} - 2u_T \right)^2 = \frac{|\omega_{K,L}^T|}{324} (5u_K - u_L - 4u_M)^2 \\ &\geq \frac{2|T|}{243} (u_K - u_M)^2 = \frac{2}{243} \int_T |\bar{u}_T(\mathbf{x}) - \underline{u}_T(\mathbf{x})|^2 \, d\mathbf{x}. \end{aligned} \quad (\text{A.4})$$

On the other hand, it follows from the Poincaré-Wirtinger inequality that (see [70])

$$\int_T |u_T(\mathbf{x}) - u_T|^2 \, d\mathbf{x} \leq \frac{h_T^2}{\pi^2} \int_T |\nabla u_T(\mathbf{x})|^2 \, d\mathbf{x},$$

which, together with (A.4), yields

$$\begin{aligned} \int_{\Omega} |\bar{u}_T(\mathbf{x}) - \underline{u}_T(\mathbf{x})|^2 \, d\mathbf{x} &= \sum_{T \in \mathcal{T}} \int_T |\bar{u}_T(\mathbf{x}) - \underline{u}_T(\mathbf{x})|^2 \, d\mathbf{x} \\ &\leq \frac{243}{2} \sum_{T \in \mathcal{T}} \int_T |u_T(\mathbf{x}) - u_T|^2 \, d\mathbf{x} \leq ch^2 \int_{\Omega} |\nabla u_T(\mathbf{x})|^2 \, d\mathbf{x}. \end{aligned}$$

This ends the proof of the lemma. \square

Remark 5. We can easily see that Lemma A.1 holds in L^1 . Indeed, using inequality (A.3) and the Poincaré-Wirtinger inequality in L^1 for convex domains (see [1]), one gets

$$\int_{\Omega} |\bar{u}_{\mathcal{T}}(\mathbf{x}) - \underline{u}_{\mathcal{T}}(\mathbf{x})| \, d\mathbf{x} \leq \frac{27h}{2} \int_{\Omega} |\nabla u_{\mathcal{T}}(\mathbf{x})| \, d\mathbf{x}.$$

Lemma A.2. Let $(u_K)_{K \in \mathcal{V}} \in \mathbb{R}^{\#\mathcal{V}}$, and let $u_{\mathcal{T}}$ and $u_{\mathcal{M}}$ be respectively the corresponding piecewise linear and piecewise constant reconstructions, then

$$\int_{\Omega} |u_{\mathcal{T}}(\mathbf{x}) - u_{\mathcal{M}}(\mathbf{x})|^2 \, d\mathbf{x} \leq ch^2 \|\nabla u_{\mathcal{T}}\|_{(L^2(\Omega))^2}^2, \quad (\text{A.5})$$

where c is the same constant given in Lemma A.1.

Proof. In order to prove inequality (A.5), we first prove the following inequality

$$|u_{\mathcal{T}}(\mathbf{x}) - u_{\mathcal{M}}(\mathbf{x})| \leq |\bar{u}_{\mathcal{T}}(\mathbf{x}) - \underline{u}_{\mathcal{T}}(\mathbf{x})|, \quad \text{for a.e. } \mathbf{x} \in \Omega. \quad (\text{A.6})$$

Let $\mathbf{x} \in \Omega = \cup_{T \in \mathcal{T}} \bar{T}$, there exists a triangle T such that $\mathbf{x} \in \bar{T}$. We assume, without loss of generality, that \mathbf{x} lies in the interior of the triangle $\omega_{K,L}^T$.

One has, from the definition of $\bar{u}_{\mathcal{T}}$ and $\underline{u}_{\mathcal{T}}$, that

$$u_{\mathcal{T}}(\mathbf{x}) - u_{\mathcal{M}}(\mathbf{x}) = u_{\mathcal{T}}(\mathbf{x}) - u_K \leq \bar{u}_{\mathcal{T}}(\mathbf{x}) - \underline{u}_{\mathcal{T}}(\mathbf{x}).$$

On the other hand, we have $u_{\mathcal{M}}(\mathbf{x}) = u_K \leq \bar{u}_{\mathcal{T}}(\mathbf{x})$ and since $\underline{u}_{\mathcal{T}}(\mathbf{x}) \leq u_{\mathcal{T}}(\mathbf{x})$, one gets

$$u_{\mathcal{M}}(\mathbf{x}) - u_{\mathcal{T}}(\mathbf{x}) \leq \bar{u}_{\mathcal{T}}(\mathbf{x}) - \underline{u}_{\mathcal{T}}(\mathbf{x}).$$

One can conclude the proof of this lemma by using Lemma A.1 as well as inequality (A.6). \square

Let $T \in \mathcal{T}$ be a triangle whose vertices are located at \mathbf{x}_K , \mathbf{x}_L , and \mathbf{x}_M . We denote by \mathbf{x}_T its center of gravity, by \mathbf{x}_{ML} the midpoint of $[\mathbf{x}_M, \mathbf{x}_L]$, and by \mathbf{x}_{MK} the midpoint of $[\mathbf{x}_M, \mathbf{x}_K]$. We denote by ω_M the dual control volume constructed around \mathbf{x}_M and by ω_M^T the intersection between ω_M and T .

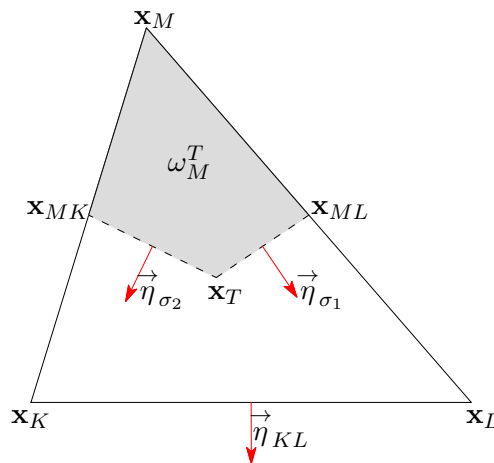


FIGURE A.2 – The triangle $T \in \mathcal{T}$ and the normal vectors on the edges.

Lemma A.3. *With the notations of figure A.2, the following geometric property*

$$\sum_{i=1}^2 |\sigma_i| \vec{\eta}_{\sigma_i} = \frac{1}{2} |\mathbf{x}_K \mathbf{x}_L| \vec{\eta}_{KL}$$

holds.

Proof. Applying the Green–Gauss formula to the constant field vector $\vec{1}$ over the domain ω_M^T , one has

$$\int_{\omega_M^T} \operatorname{div} \left(\vec{1} \right) d\mathbf{x} = 0 = \int_{\partial\omega_M^T} \vec{\eta} d\sigma = \sum_{\sigma \subset \partial\omega_M^T} |\sigma| \vec{\eta}_\sigma = \sum_{i=1}^4 |\sigma_i| \vec{\eta}_{\sigma_i}, \quad (\text{A.7})$$

where σ_3 (resp. σ_4) represents the length of $[\mathbf{x}_M, \mathbf{x}_{ML}]$ (resp. $[\mathbf{x}_M, \mathbf{x}_{MK}]$) and $\vec{\eta}_{\sigma_i}$ represents the unit normal vector to σ_i , $i = 3, 4$ outward to T .

On the other hand, one applies again the Green–Gauss formula over the triangle T and gets

$$\begin{aligned} \int_T \operatorname{div} \left(\vec{1} \right) d\mathbf{x} &= 0 = \int_{\partial T} \vec{\eta} d\sigma = \sum_{\sigma \subset \partial T} |\sigma| \vec{\eta}_\sigma \\ &= |\mathbf{x}_M \mathbf{x}_L| \vec{\eta}_{\sigma_3} + |\mathbf{x}_M \mathbf{x}_K| \vec{\eta}_{\sigma_4} + |\mathbf{x}_K \mathbf{x}_L| \vec{\eta}_{KL}. \end{aligned} \quad (\text{A.8})$$

One can conclude the proof of the lemma by combining equations (A.7)–(A.8). \square

A.1 The reference element

We recall that in *the control volume finite element method*, we have introduced the transmissibility coefficients

$$\Lambda_{KL} = \int_{\Omega} \Lambda(\mathbf{x}) \nabla \varphi_K(\mathbf{x}) \cdot \nabla \varphi_L(\mathbf{x}) d\mathbf{x}.$$

In a computation point of view, and in order to approximate these elements, we need : either (a) an effective way of evaluating the functions $(\varphi_K)_{K \in \mathcal{V}}$ and their gradients, or (b) a closed form for the resulting integrals. Both possibilities are done usually by introducing the so-called reference element.

For triangles, the reference element \hat{T} is the triangle with vertices

$$\widehat{S}_1 = (0, 0), \quad \widehat{S}_2 = (1, 0), \quad \widehat{S}_3 = (0, 1).$$

To distinguish variables in the reference element and in a physical triangle (a general triangle), we usually adopt the variables $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ in the reference element and (\mathbf{x}, \mathbf{y}) in the physical element. We designate by T a physical triangle with vertices $S_1 = (\mathbf{x}_1, \mathbf{y}_1)$, $S_2 = (\mathbf{x}_2, \mathbf{y}_2)$, and $S_3 = (\mathbf{x}_3, \mathbf{y}_3)$. We denote by $(\mathbf{x}_{ji}, \mathbf{y}_{ji}) = (\mathbf{x}_j - \mathbf{x}_i, \mathbf{y}_j - \mathbf{y}_i)$ the coordinates of the vector $\vec{S_i S_j}$, for every $i \neq j \in \{1, 2, 3\}$.

Let \mathcal{F}_T be the affine transformation mapping the triangle \hat{T} bijectively into T . We have $\mathcal{F}_T(\widehat{S}_i) = S_i$ for every $i = 1, 2, 3$. Thus we have

$$\begin{aligned} \mathcal{F}_T \begin{pmatrix} \hat{\mathbf{x}} \\ \hat{\mathbf{y}} \end{pmatrix} &= \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_{21} & \mathbf{x}_{31} \\ \mathbf{y}_{21} & \mathbf{y}_{31} \end{pmatrix} \begin{pmatrix} \hat{\mathbf{x}} \\ \hat{\mathbf{y}} \end{pmatrix} + \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{y}_1 \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{y}_1 \end{pmatrix} (1 - \hat{\mathbf{x}} - \hat{\mathbf{y}}) + \begin{pmatrix} \mathbf{x}_2 \\ \mathbf{y}_2 \end{pmatrix} \hat{\mathbf{x}} + \begin{pmatrix} \mathbf{x}_3 \\ \mathbf{y}_3 \end{pmatrix} \hat{\mathbf{y}}. \end{aligned} \quad (\text{A.9})$$

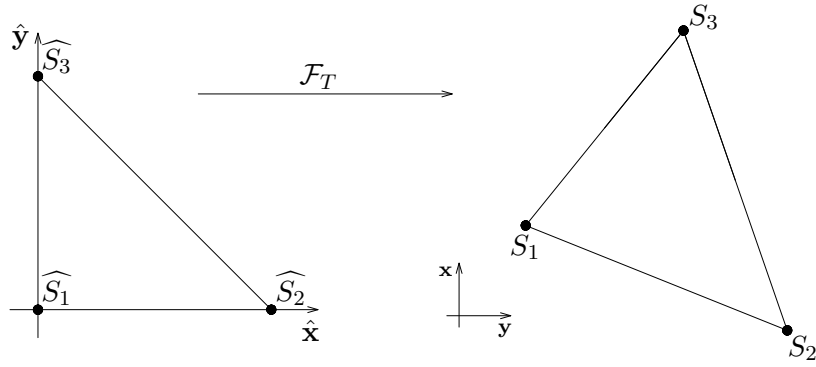


FIGURE A.3 – Reference and physical triangles.

The local nodal functions in the reference triangles are three \mathbb{P}_1 functions satisfying

$$\widehat{\varphi}_i(\widehat{S}_j) = \delta_{ij}, \quad i, j = 1, 2, 3.$$

They are precisely defined, into the space of coordinates $(\widehat{x}, \widehat{y})$, by

$$\widehat{\varphi}_1(\widehat{x}, \widehat{y}) = 1 - \widehat{x} - \widehat{y},$$

$$\widehat{\varphi}_2(\widehat{x}, \widehat{y}) = \widehat{x},$$

$$\widehat{\varphi}_3(\widehat{x}, \widehat{y}) = \widehat{y}.$$

It is simple now to prove that

$$\widehat{\varphi}_i = \varphi_i \circ \mathcal{F}_T, \quad i = 1, 2, 3.$$

or, what is the same

$$\varphi_i = \widehat{\varphi}_i \circ \mathcal{F}_T^{-1}, \quad i = 1, 2, 3. \quad (\text{A.10})$$

We denote by \mathbf{B}_T the matrix of the linear transformation \mathcal{F}_T , we have

$$\det(\mathbf{B}_T) = 2|T| \neq 0 \quad \text{and} \quad \mathcal{F}_T^{-1} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \begin{pmatrix} \widehat{x} \\ \widehat{y} \end{pmatrix} = \mathbf{B}_T^{-1} \begin{pmatrix} \mathbf{x} - \mathbf{x}_1 \\ \mathbf{y} - \mathbf{y}_1 \end{pmatrix}, \quad (\text{A.11})$$

where $\mathbf{B}_T^{-1} = \frac{1}{2|T|} \begin{pmatrix} \mathbf{y}_{31} & -\mathbf{x}_{31} \\ -\mathbf{y}_{21} & \mathbf{x}_{21} \end{pmatrix}$ is the inverse of the matrix \mathbf{B}_T .

Using together equations (A.10) and (A.11), one has a simple way of evaluating the nodal functions $(\varphi_i)_{i=1,2,3}$ and therefore their gradients. For instance, one gets

$$\varphi_1(\mathbf{x}, \mathbf{y}) = 1 - \frac{1}{2|T|} (\mathbf{y}_{32}(\mathbf{x} - \mathbf{x}_1) - \mathbf{x}_{32}(\mathbf{y} - \mathbf{y}_1)),$$

$$\varphi_2(\mathbf{x}, \mathbf{y}) = \frac{1}{2|T|} (\mathbf{y}_{31}(\mathbf{x} - \mathbf{x}_1) - \mathbf{x}_{31}(\mathbf{y} - \mathbf{y}_1)),$$

$$\varphi_3(\mathbf{x}, \mathbf{y}) = \frac{1}{2|T|} (-\mathbf{y}_{21}(\mathbf{x} - \mathbf{x}_1) + \mathbf{x}_{21}(\mathbf{y} - \mathbf{y}_1)),$$

and

$$\nabla\varphi_1(\mathbf{x}, \mathbf{y}) = \frac{-1}{2|T|} \begin{pmatrix} \mathbf{y}_{32} \\ -\mathbf{x}_{32} \end{pmatrix} = \frac{-\|\vec{S_2 S_3}\|}{2|T|} \vec{\eta}_{23},$$

$$\nabla\varphi_2(\mathbf{x}, \mathbf{y}) = \frac{1}{2|T|} \begin{pmatrix} \mathbf{y}_{31} \\ -\mathbf{x}_{31} \end{pmatrix} = \frac{-\|\vec{S_3 S_1}\|}{2|T|} \vec{\eta}_{31},$$

$$\nabla\varphi_3(\mathbf{x}, \mathbf{y}) = \frac{-1}{2|T|} \begin{pmatrix} \mathbf{y}_{21} \\ -\mathbf{x}_{21} \end{pmatrix} = \frac{-\|\vec{S_1 S_2}\|}{2|T|} \vec{\eta}_{12},$$

where $\vec{\eta}_{ij} = (-\mathbf{y}_{ij}, \mathbf{x}_{ij})$ is the unit normal vector at the line segment $[S_i S_j]$ **outward** to T , for every $i \neq j \in \{1, 2, 3\}$.

Bibliographie

- [1] Gabriel Acosta and Ricardo G Durán. An optimal poincaré inequality in L^1 for convex domains. *Proceedings of the American Mathematical Society*, pages 195–202, 2004. [153](#)
- [2] Mohamed Afif and Brahim Amaziane. Convergence of finite volume schemes for a degenerate convection–diffusion equation arising in flow in porous media. *Computer methods in applied mechanics and engineering*, 191(46) :5265–5286, 2002. [54](#), [55](#)
- [3] Abdellatif Agouzal, Jacques Baranger, Jean-François Maitre, and Fabienne Oudin. Connection between finite volume and mixed finite element methods for a diffusion problem with nonconstant coefficients. application to a convection diffusion problem. *EAST WEST J NUMER MATH*, 3(4) :237–254, 1995. [54](#)
- [4] Herbert Amann. Nonhomogeneous linear and quasilinear elliptic and parabolic boundary value problems. *Function spaces, differential operators and nonlinear analysis*, 133 :9–126, 1993. [32](#)
- [5] Youcef Amirat, Kamel Hamdache, and Abdelhamid Ziani. Mathematical analysis for compressible miscible displacement models in porous media. *Mathematical Models and Methods in Applied Sciences*, 6(06) :729–747, 1996. [120](#)
- [6] Youcef Amirat and Abdelhamid Ziani. Global weak solutions for a parabolic system modeling a one-dimensional miscible flow in porous media. *Journal of mathematical analysis and applications*, 220(2) :697–718, 1998. [120](#)
- [7] Boris Andreianov, Mostafa Bendahmane, and Mazen Saad. Finite volume methods for degenerate chemotaxis model. *Journal of computational and applied mathematics*, 235(14) :4015–4031, 2011. [31](#), [47](#), [54](#), [55](#)
- [8] Mostafa Bendahmane, Kenneth H Karlsen, and José Miguel Urbano. On a two-sidedly degenerate chemotaxis model with volume-filling effect. *Mathematical Models and Methods in Applied Sciences*, 17(05) :783–804, 2007. [4](#), [31](#), [120](#)
- [9] Mostafa Bendahmane and Mazen Saad. Mathematical analysis and pattern formation for a partial immune system modeling the spread of an epidemic disease. *Acta applicandae mathematicae*, 115(1) :17–42, 2011. [30](#), [54](#)
- [10] Fred Brauer and Carlos Castillo-Chavez. *Mathematical models in population biology and epidemiology*. Springer, 2011. [56](#)
- [11] Konstantin Brenner, Roland Masson, et al. Convergence of a vertex centred discretization of two-phase darcy flows on general meshes. *International Journal of Finite Volume*, 10 :1–37, 2013. [96](#)
- [12] Haïm Brezis. *Analyse fonctionnelle : Théorie et applications*, volume 5. Masson, 1983. [74](#)

- [13] Giuseppe Buttazzo, Gérard Michaille, and Hedy Attouch. *Variational analysis in Sobolev and BV spaces : applications to PDEs and optimization*, volume 6. Siam, 2006. [74](#)
- [14] HM Byrne and MAJ Chaplain. Mathematical models for tumour angiogenesis : numerical simulations and nonlinear wave solutions. *Bulletin of Mathematical Biology*, 57(3) :461–486, 1995. [30](#), [54](#)
- [15] Zhiqiang Cai. On the finite volume element method. *Numerische Mathematik*, 58(1) :713–735, 1990. [54](#)
- [16] Clément Cancès and Cindy Guichard. Convergence of a nonlinear entropy diminishing control volume finite element scheme for solving anisotropic degenerate parabolic equations. HAL : hal-00955091, 2014. [73](#), [86](#), [94](#), [103](#), [104](#), [105](#)
- [17] Emilio Cariaga, Fernando Concha, Iuliu Sorin Pop, and Mauricio Sepúlveda. Convergence analysis of a vertex-centered finite volume scheme for a copper heap leaching model. *Mathematical Methods in the Applied Sciences*, 33(9) :1059–1077, 2010. [54](#), [55](#)
- [18] Juan Casado-Díaz, T Chacón Rebollo, Vivette Girault, M Gómez Mármol, and François Murat. Finite elements approximation of second order linear elliptic equations in divergence form with right-hand side in L^1 . *Numerische Mathematik*, 105(3) :337–374, 2007. [63](#)
- [19] Mark AJ Chaplain, Mahadevan Ganesh, and Ivan G Graham. Spatio-temporal pattern formation on spherical surfaces : numerical simulation and application to solid tumour growth. *Journal of mathematical biology*, 42(5) :387–423, 2001. [30](#), [54](#)
- [20] Guy Chavent, Jérôme Jaffré, and Jean E Roberts. Mixed-hybrid finite elements and cell-centred finite volumes for two-phase flow in porous media. *Mathematical Modelling of Flow Through Porous Media*, pages 100–114, 1995. [54](#)
- [21] John Condeelis, Robert H Singer, and Jeffrey E Segall. The great escape : when cancer cells hijack the genes for chemotaxis and motility. *Annu. Rev. Cell Dev. Biol.*, 21 :695–718, 2005. [4](#)
- [22] Yves Coudière, Jean-Paul Vila, and Philippe Villedieu. Convergence rate of a finite volume scheme for a two dimensional convection-diffusion problem. *ESAIM : Mathematical Modelling and Numerical Analysis*, 33(03) :493–516, 1999. [54](#)
- [23] Klaus Deimling. Nonlinear functional analysis. 1985. *Springer-Verlag, Berlin*. [105](#)
- [24] Dirk Dormann and Cornelis J Weijer. Chemotactic cell movement during development. *Current opinion in genetics & development*, 13(4) :358–364, 2003. [4](#)
- [25] Jérôme Droniou, Robert Eymard, Thierry Gallouët, and Raphaelle Herbin. Gradient schemes : a generic framework for the discretisation of linear, nonlinear and nonlocal elliptic and parabolic equations. *Mathematical Models and Methods in Applied Sciences*, 23(13) :2395–2432, 2013. [54](#)
- [26] Renjun Duan, Alexander Lorz, and Peter Markowich. Global solutions to the coupled chemotaxis-fluid equations. *Communications in Partial Differential Equations*, 35(9) :1635–1673, 2010. [120](#)
- [27] Alexandre Ern and Jean-Luc Guermond. Theory and practice of finite elements, vol. 159 of applied mathematical sciences, 2004. [96](#)

- [28] Lawrence C Evans. Partial differential equations : Graduate studies in mathematics. *American Mathematical Society*, 2, 1998. 69, 125
- [29] Robert Eymard and Thierry Gallouët. Convergence d'un schéma de type éléments finis-volumes finis pour un système formé d'une équation elliptique et d'une équation hyperbolique. *Modélisation mathématique et analyse numérique*, 27(7) :843–861, 1993. 54
- [30] Robert Eymard, Thierry Gallouët, Mustapha Ghilani, and Raphaële Herbin. Error estimates for the approximate solutions of a nonlinear hyperbolic equation given by finite volume schemes. *IMA J. Numer. Anal.*, 18(4) :563–594, 1998. 105
- [31] Robert Eymard, Thierry Gallouët, and Raphaële Herbin. Finite volume methods. *Handbook of numerical analysis*, 7 :713–1018, 2000. 7, 41, 46, 54, 75, 91
- [32] Robert Eymard, Thierry Gallouët, and Raphael Herbin. A finite volume scheme for anisotropic diffusion problems. *Comptes Rendus Mathématique*, 339(4) :299–302, 2004. 54, 92, 109
- [33] Robert Eymard, Thierry Gallouët, and Raphael Herbin. Discretization of heterogeneous and anisotropic diffusion problems on general nonconforming meshes sushi : a scheme using stabilization and hybrid interfaces. *IMA Journal of Numerical Analysis*, 30(4) :1009–1043, 2010. 54
- [34] Robert Eymard, Thierry Gallouët, Raphaële Herbin, and Anthony Michel. Convergence of a finite volume scheme for nonlinear degenerate parabolic equations. *Numerische Mathematik*, 92(1) :41–82, 2002. 110
- [35] Pierre Fabrie and Michel Langlais. Mathematical analysis of miscible displacement in porous medium. *SIAM journal on mathematical analysis*, 23(6) :1375–1392, 1992. 120
- [36] Francis Filbet. A finite volume scheme for the patlak–keller–segel chemotaxis model. *Numerische Mathematik*, 104(4) :457–488, 2006. 47
- [37] Peter A Forsyth. A control volume finite element approach to napl groundwater contamination. *SIAM Journal on Scientific and Statistical Computing*, 12(5) :1029–1057, 1991. 55
- [38] Thierry Gallouët and Jean-Claude Latché. Compactness of discrete approximate solutions to parabolic pdes-application to a turbulence model. *Communications on Pure & Applied Analysis*, 11(6), 2012. 109
- [39] Cédric Galusinski and Mazen Saad. On a degenerate parabolic system for compressible, immiscible, two-phase flows in porous media. *Advances in Differential Equations*, 9(11-12) :1235–1278, 2004. 120
- [40] Sergueï Godounov, A Zabrodine, Mikhail Ivanov, A Kraiko, and G Prokopov. *Résolution numérique des problèmes multidimensionnels de la dynamique des gaz*. Editions Mir, 1979. 92
- [41] SK Godunov, A Zabrodine, M Ivanov, A Kraiko, and G Prokopov. *Résolution numérique des problèmes multidimensionnels de la dynamique des gaz*. 1979. *Mir, Moscow*. 91
- [42] Raphaële Herbin and Florence Hubert. Benchmark on discretization schemes for anisotropic diffusion problems on general grids. *Finite volumes for complex applications V*, pages 659–692, 2008. 55

- [43] E Hildebrand and UB Kaupp. Sperm chemotaxis : a primer. *Annals of the New York Academy of Sciences*, 1061(1) :221–225, 2005. [4](#)
- [44] Thomas Hillen and Kevin Painter. Global existence for a parabolic chemotaxis model with prevention of overcrowding. *Advances in Applied Mathematics*, 26(4) :280–301, 2001. [5](#)
- [45] Thomas Hillen and Kevin J Painter. A user’s guide to pde models for chemotaxis. *Journal of mathematical biology*, 58(1-2) :183–217, 2009. [56](#)
- [46] Dirk Horstmann. From 1970 until present : the keller-segel model in chemotaxis and its consequences. i. *Jahresbericht der Deutschen Mathematiker Vereinigung*, 105 :103–165, 2003. [31](#)
- [47] Dirk Horstmann. From 1970 until present : The keller-segel model in chemotaxis and its consequences ii. *Jahresbericht der Deutschen Mathematiker Vereinigung*, 106(2) :51–70, 2004. [4](#), [30](#), [31](#), [54](#)
- [48] Hyung Ju Hwang, Kyungkeun Kang, and Angela Stevens. Drift-diffusion limits of kinetic models for chemotaxis : a generalization. *Discrete Contin. Dyn. Syst. Ser. B*, 5(2) :319–334, 2005. [37](#)
- [49] Moustafa Ibrahim and Mazen Saad. Pattern formation and cross-diffusion for a chemotaxis model. In *Conference Book Biomath 2013*, 2013. [79](#)
- [50] Moustafa Ibrahim and Mazen Saad. On the efficacy of a control volume finite element method for the capture of patterns for a volume-filling chemotaxis model. *Computers & Mathematics with Applications*, 2014. [86](#), [93](#)
- [51] Evelyn F Keller and Lee A Segel. Initiation of slime mold aggregation viewed as an instability. *Journal of Theoretical Biology*, 26(3) :399–415, 1970. [4](#), [30](#), [54](#)
- [52] Evelyn F Keller and Lee A Segel. Model for chemotaxis. *Journal of Theoretical Biology*, 30(2) :225–234, 1971. [30](#), [54](#)
- [53] Bruno Larrivee and Aly Karsan. Signaling pathways induced by vascular endothelial growth factor (review). *International journal of molecular medicine*, 5(5) :447–503, 2000. [4](#)
- [54] Jean Leray and Juliusz Schauder. Topologie et équations fonctionnelles. *Ann. Sci. École Norm. Sup.(3)*, 51 :45–78, 1934. [105](#)
- [55] Jacques-Louis Lions, Jacques-Louis Lions, Jacques-Louis Lions, and Jacques-Louis Lions. *Quelques méthodes de résolution des problèmes aux limites non linéaires*, volume 31. Dunod Paris, 1969. [69](#)
- [56] Pierre-Louis Lions. Mathematical topics in fluid mechanics : Volume 1 : Incompressible models. 1996. [125](#)
- [57] Alexander Lorz. Coupled chemotaxis fluid model. *Mathematical Models and Methods in Applied Sciences*, 20(06) :987–1004, 2010. [120](#)
- [58] Philip K Maini. Using mathematical models to help understand biological pattern formation. *Comptes rendus biologiques*, 327(3) :225–234, 2004. [30](#)
- [59] James D Murray. How the leopard gets its spots. *Scientific American*, 258(3) :80–87, 1988. [2](#), [54](#)

- [60] James D Murray. *Mathematical biology : I. An introduction*, volume 17 of *Interdisciplinary Applied Mathematics*. Springer, 2002. [30](#), [33](#)
- [61] James D Murray. *Mathematical biology II : spatial models and biomedical applications*, volume 18 of *Interdisciplinary Applied Mathematics*. Springer, New York, 2003. [2](#), [3](#), [30](#), [33](#), [34](#)
- [62] James D Murray. Turing centennial celebration. Princeton University, May 10-12 2012. [2](#)
- [63] James D Murray and MR Myerscough. Pigmentation pattern formation on snakes. *Journal of theoretical biology*, 149(3) :339–360, 1991. [30](#), [54](#)
- [64] Mario Ohlberger. A posteriori error estimates for vertex centered finite volume approximations of convection-diffusion-reaction equations. *ESAIM : Mathematical Modelling and Numerical Analysis*, 35(02) :355–387, 2001. [55](#)
- [65] Stanley Osher and Fred Solomon. Upwind difference schemes for hyperbolic systems of conservation laws. *Mathematics of computation*, 38(158) :339–374, 1982. [92](#)
- [66] Kevin J Painter and Thomas Hillen. Volume-filling and quorum-sensing in models for chemosensitive movement. *Can. Appl. Math. Quart.*, 10(4) :501–543, 2002. [31](#), [56](#)
- [67] KJ Painter, PK Maini, and HG Othmer. Stripe formation in juvenile pomacanthus explained by a generalized turing mechanism with chemotaxis. *Proceedings of the National Academy of Sciences*, 96(10) :5549–5554, 1999. [30](#), [54](#)
- [68] Hwan Tae Park, Jane Wu, and Yi Rao. Molecular control of neuronal migration. *Bioessays*, 24(9) :821–827, 2002. [4](#)
- [69] Clifford S Patlak. Random walk with persistence and external bias. *The Bulletin of mathematical biophysics*, 15(3) :311–338, 1953. [4](#), [30](#), [54](#)
- [70] Lawrence E Payne and Hans F Weinberger. An optimal poincaré inequality for convex domains. *Archive for Rational Mechanics and Analysis*, 5(1) :286–292, 1960. [152](#)
- [71] Alex Potapov and Thomas Hillen. Metastability in chemotaxis models. *Journal of Dynamics and Differential Equations*, 17(2) :293–330, 2005. [31](#)
- [72] P Schaaf and J Talbot. Surface exclusion effects in adsorption processes. *The Journal of chemical physics*, 91 :4401, 1989. [31](#)
- [73] Jacques Simon. Compact sets in the space $L^p(0, T; B)$. *Annali di Matematica pura ed applicata*, 146(1) :65–96, 1986. [124](#)
- [74] Angela Stevens and Hans G Othmer. Aggregation, blowup, and collapse : The abc’s of taxis in reinforced random walks. *SIAM Journal on Applied Mathematics*, 57(4) :1044–1081, 1997. [56](#)
- [75] Roger Temam. Navier-stokes equations. theory and numerical analysis. reprint of the 1984 edition. *AMS Chelsea, Providence, RI*, 2001. [69](#)
- [76] Alan Mathison Turing. The chemical basis of morphogenesis. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, 237(641) :37–72, 1952. [2](#), [30](#), [54](#)

- [77] R Tyson, SR Lubkin, and James D Murray. A minimal mechanism for bacterial pattern formation. *Proceedings of the Royal Society of London. Series B : Biological Sciences*, 266(1416) :299–304, 1999. [53](#)
- [78] Zhian Wang and Thomas Hillen. Classical solutions and pattern formation for a volume filling chemotaxis model. *Chaos : An Interdisciplinary Journal of Nonlinear Science*, 17(3) :037108–037108, 2007. [5](#), [7](#), [31](#), [32](#), [56](#), [79](#)
- [79] Dariusz Wrzosek. Global attractor for a chemotaxis model with prevention of overcrowding. *Nonlinear Analysis : Theory, Methods & Applications*, 59(8) :1293–1310, 2004. [5](#)
- [80] Xuesong Yang, Dirk Dormann, Andrea E Münsterberg, and Cornelis J Weijer. Cell movement patterns during gastrulation in the chick are controlled by positive and negative chemotaxis mediated by fgf4 and fgf8. *Developmental cell*, 3(3) :425–437, 2002. [30](#)
- [81] Anatol M Zhabotinsky. Periodical oxidation of malonic acid in solution (a study of the belousov reaction kinetics). *Biofizika*, 9 :306–311, 1964. [29](#)

Thèse de Doctorat

Moustafa IBRAHIM

Systèmes paraboliques dégénérés intervenant en mécanique des fluides et en médecine: analyse mathématique et numérique

Degenerate parabolic systems involved in fluid mechanics and medicine: mathematical and numerical analysis

Résumé

Dans cette thèse, nous nous intéressons à l'analyse mathématique et numérique des systèmes paraboliques non linéaires dégénérés découlant, soit de la modélisation de la chimiotaxie, soit de la modélisation des fluides compressibles. Le modèle de chimiotaxie (Keller-Segel) proposé est un modèle de dynamique des populations décrivant l'évolution spatio-temporelle de la densité cellulaire et de la concentration chimiotactique. Pour ce modèle, nous étudions la formation de patterns en utilisant l'analyse de stabilité linéaire et le principe de Turing. Nous proposons ensuite un schéma numérique CVFE pour un modèle anisotrope de Keller-Segel. La construction de ce schéma est basée sur la méthode des éléments finis pour le terme de diffusion et sur la méthode des volumes finis classique pour le terme de convection. Nous montrons que ce schéma assure le principe de maximum discret et qu'il est consistant dans le cas où tous les coefficients de transmissibilité sont positifs. Par la suite, sur des maillages triangulaires généraux, nous proposons et analysons un schéma numérique CVFE non linéaire. Ce schéma est basé sur l'utilisation d'un flux numérique de Godunov pour le terme de diffusion, tandis que le terme de convection est approché au moyen d'un décentrage amont et d'un flux de Godunov. D'une part, le décentrage amont permet d'avoir le principe de maximum. D'autre part, le flux de Godunov assure que les solutions discrètes soient bornées sans restriction sur le maillage du domaine spatial ni sur les coefficients de transmissibilité. Nous réalisons différentes simulations numériques bi-dimensionnelles pour illustrer l'efficacité du schéma à tenir compte des hétérogénéités. Enfin, nous nous intéressons à une équation parabolique dégénérée contenant des termes dégénérés d'ordre 0 et 1 et décrivant un modèle de chimiotaxie-fluide ou l'écoulement d'un fluide compressible. Une formulation faible classique est souvent possible en absence des termes dégénérés d'ordre 0 et 1 ; tandis que dans le cas général, nous obtenons des solutions dans un sens affaibli vérifiant une formulation de type inégalité variationnelle. La définition des solutions faibles est adaptée à la nature de la dégénérescence des termes de dissipation.

Mots clés

Systèmes de réaction-diffusion, Formation de patterns, Chimiotaxie, Stabilité linéaire, Méthode de volumes finis, Méthode des éléments finis, Systèmes paraboliques dégénérés, Tenseurs anisotropes hétérogènes, Fluide compressible.

Abstract

In this thesis, we are interested in the mathematical and numerical analysis of nonlinear degenerate parabolic systems arising either from modeling the chemotaxis process, or from modeling compressible flows in porous media. The proposed chemotaxis model (Keller-Segel model) is a model of population dynamics describing the spatio-temporal evolution of the cell density and the chemical concentration. For this model, we study the pattern formation using the linear stability analysis as well as the principle of Turing. Then, we propose a numerical scheme (CVFE scheme) for an anisotropic Keller-Segel model. The construction of the scheme is based on the use of each of the finite element scheme for the diffusion term and the upwind finite volume scheme for the convective term. We show that the scheme is consistent and ensures the discrete maximum principle in the case where all the transmissibility coefficients are nonnegative. Thereafter, over general triangular meshes, we propose and analyze a nonlinear CVFE scheme. This scheme is based on the use of the Godunov flux function for the diffusion term, while the convective term is approximated by parts using an upwind finite volume scheme and a Godunov flux function. First, the upwind finite volume scheme allows of having the discrete maximum principle. On the other hand, the Godunov scheme ensures the boundedness of the discrete solutions without restrictions on the mesh nor on the transmissibility coefficients. Using this scheme, we realize some numerical simulations to illustrate the effectiveness of the scheme. Finally, we are interested in a degenerate parabolic equation containing degenerate terms of order 0 and 1 and describing a chemotaxis-fluid model or a displacement of compressible flows. Classical weak formulation is often possible in the absence of degenerate terms of order 0 and 1; while in the general case, we obtain weak solutions in the sense of verifying a weighted formulation. The definition of weak solutions is adapted to the nature of the degeneracy of the dissipative terms.

Key Words

Reaction-diffusion systems, Pattern formation, Chemotaxis, Linear stability, Finite volume method, Finite element method, Degenerate parabolic systems, Heterogeneous anisotropic tensors, Compressible fluid.